

HOP: An Efficient Homotopy Method for Online Anisotropic Projection

Yuting Ye *

Division of Biostatistics
Univ. of California, Berkeley

Lihua Lei †

Department of Statistics
Univ. of California, Berkeley

Cheng Ju

Division of Biostatistics
Univ. of California, Berkeley

Abstract

In this article, we develop and analyze a homotopy continuation method for a generic online quadratic programming problem, which has several important applications such as Markowitz portfolio selection and Online Newton Step. We refer to our method as *Homotopic Online Projection* (HOP) algorithm. HOP always produces exact solutions and is computationally efficient especially when the solutions are sparse or the solutions change slowly over the time. The superior performance is shown on both synthetic datasets and real datasets.

1 Introduction

The problem of minimizing a convex quadratic objective subjected to linear constraints has been studied for decades, e.g.[17, 25, 28]. This type of problem has a generic form as

$$\min \frac{1}{2}x^T A x - r^T x, \quad s.t. \quad Bx = b, x \geq 0, \quad (1)$$

where A is a positive semi-definite matrix in $\mathbb{R}^{n \times n}$, $x, y \in \mathbb{R}^n$, $B \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$. The notion $x \geq 0$ means that all elements of x are non-negative. One case of particular importance involves $B = \mathbf{1}^T$ where $\mathbf{1}$ is a n -dimensional vector with all entries 1. This constraint amounts to restricting the solution into the $(n - 1)$ -dimensional simplex. In this case, x can be interpreted as a probability distribution or an assignment weight, both of which are natural objects in many applications [24, 21, 12, 27, 9, 6].

From a geometric view, (1) is equivalent to finding the projection of $y = A^{-1}r$ onto the simplex with an anisotropic norm $\|\cdot\|_A$ such that $\|z\|_A = \sqrt{z^T A z}$. In the isotropic case where $A = I$, the projection is done in the Euclidean space and the problem can be solved quite efficiently in $O(n)$ time [13] due to the explicit form of the solution. For general anisotropic matrix A , such benefits disappear and one must resort to iterative algorithms, such as interior-point method [29]. Although not as fast as the isotropic case, those algorithms are still effective and stable.

In recent years, online problems attract a lot of interests. In many online problems, one need to solve a sequence of quadratic programming problems in (1), formulated as

$$\min_{x \in \Delta_n} \frac{1}{2}x^T A^{(t)} x - (r^{(t)})^T x \quad (2)$$

where $\{A^{(t)} \in \mathbb{R}^{n \times n} : t = 1, 2, \dots\}$ is a stream of positive semi-definitive matrices, $\{r^{(t)} \in \mathbb{R}^n : t = 1, 2, \dots\}$ is a stream of vectors, and Δ_n is the $(n - 1)$ -dimensional simplex. A naive approach would be to apply any off-the-shelf offline algorithm to solve (2) sequentially. However, it could be inefficient when $A^{(t)}$ and $r^{(t)}$ have some stable structures along with the time. In this article, we consider a particular case where $A^{(t)}$ is updated by rank-one matrix in each time, i.e.

$$A^{(t+1)} = A^{(t)} + g^{(t)}(g^{(t)})^T. \quad (3)$$

*yeyt@berkeley.edu

†lihua.lei@berkeley.edu

The problem (2) with matrix update (3) forms a basis in various applications. For instance, the mean-Variance analysis, also known as Markowitz portfolio problem, is proposed by Markowitz in his pioneering work [24] to formalize the fundamental tradeoff between the return and the risk. In the area of risk management where the risk is of primary concern, a trader aims to select a portfolio from a set of assets with low variability while maintaining a high expected return. More precisely, one is given a basket of n assets and the objective can be formulated as

$$\min_{x \in \Delta_n} \frac{1}{2} x^T \hat{\Sigma}^{(t)} x - \eta (\hat{\mu}^{(t)})^T x, \quad (4)$$

where $\hat{\Sigma}^{(t)}$ and $\hat{\mu}^{(t)}$ are the estimates of variability and average return at time t and η is a positive term characterizing the level of penalty. We will discuss the estimate of $\hat{\Sigma}^{(t)}$ and $\hat{\mu}^{(t)}$ in section 3 and show that it satisfies (3) after rescaling. Another important example is Online Newton Step (ONS), proposed by [20], which produces the iterate x_t as

$$x_t = \arg \min_{x \in \mathcal{K}} \frac{1}{2} x^T A^{(t)} x - (A^{(t)} x_{t-1} - \eta g^{(t)})^T x. \quad (5)$$

where $g^{(t)}$ is the gradient evaluated at step $t-1$, η is the stepsize and $A^{(t)} = A^{(t-1)} + g^{(t)}(g^{(t)})^T$; see [19] for details.

Heuristically, $A^{(t)}$ is perturbed in only one direction in each time and the optimal solutions in consecutive steps should be close in some sense. A usual strategy to take advantage of the minor change is the warm-start, i.e. initializing the iterate as the optimal solution in the last step. The spectral projected gradient (SPG) method [5], the projected quasi-Newton (PQN) method [23] are two state-of-the-art algorithms falling into this category. However, warm-start is not always a feasible approach. For methods like the interior-point method [29], the initial point is required to be an interior point of the constraint set, but as shown in various settings and applications, including our experiments in section 3, the solution in each step often lies on the boundary of the simplex. To the best of our knowledge, no existing method is tailored for the matrix flow (3).

To make maximal use of the structure (3), we resort to homotopy method, which is proposed decades ago and widely used in optimizing highly non-convex problems such as polynomial systems [11, 22]. The basic idea is to construct a bivariate function $H(x, w)$ on $\mathbb{R}^n \times [0, 1]$ with $H(x, 0) = g(x)$ and $H(x, 1) = f(x)$ so that in order to optimize $f(x)$ one can start from the optimizer of $g(x)$ and move towards $f(x)$ by gradually increasing w . Given sufficient smoothness of H , one can obtain a smooth trajectory, referred to as solution path in literature, penetrating the optimizers of $H(x, w)$ for all $w \in [0, 1]$ with the optimizer of $f(x)$ being the ending point. Recent years have witnessed success of homotopy methods in solving statistical problems like penalized least squares. Among others, LARS algorithm proposed by [14] to solve LASSO regression is essentially a homotopy method. Instead of the simplex, LASSO objective is a quadratic programming problem with constraint into the L_1 ball. When the data is arriving in a sequential manner, the corresponding problem, termed as online LASSO problem by [18], can be efficiently solved by constructing a homotopy between the consecutive steps.

For the problem (2) with matrix flow (3), we define the homotopy functions $H_1^{(t)}(\lambda)$ and $H_2^{(t)}(\lambda)$ as in section 2.2.1. Then solving (2) at step t amounts to optimizing $H_1^{(t)}(0)$ while solving (2) at step $t+1$ amounts to optimizing $H_2^{(t)}(1)$. In addition, notice that $H_1^{(t)}(1) = H_2^{(t)}(0)$, by optimizing $H_1^{(t)}(\lambda)$ and $H_2^{(t)}(\lambda)$ for all $\lambda, \underline{\lambda} \in [0, 1]$, we finally move from the solution at step t to that at step $t+1$. In section 2, we will show that both trajectories can be calculated efficiently, with only machine error, by utilizing the Karush-Kuhn-Tucker (KKT) conditions. More precisely, the trajectories are proved to be piecewise linear and the computation complexity for the implementation between two consecutive steps is $O(nsk)$, where s is the sparsity of the solution and k is the number of turning points, which is usually small in practice. In contrast, the rate of SPG and PQN is $O(n^2T)$ and the rate of the interior-point method is $O(n^3T)$ where T is the number of iterations. Generally speaking, the theoretical rates are incomparable but it is clear that our HOP algorithm is adaptive to the solution sparsity. In addition, we will show the exact constant hidden in the big-O notation for the HOP algorithm in section 2.3 while the constants for other algorithms are unknown and could be potentially large due to the line search step (involved in SPG and PQN). For this reason, we suggest comparing the running time on real datasets for informative comparison. Moreover, unlike the iterative methods

which can only find an approximated solution, our HOP algorithm is able to find the exact solution. In other words, the HOP algorithm can solve the problem efficiently and accurately.

The rest of the paper is organized as follows: in section 2, we propose the general framework of our HOP algorithm (Homotopic Online Projection) followed by a theoretical justification and a complexity analysis. To improve the efficiency, the real implementation is much more complicated than the general idea and hence stated in Appendices. In section 3, we show the experimental results for both synthetic data and real data. The section 4 concludes the article and discusses several possible extensions.

2 Proposed Algorithm

A generic framework to solve problem (2) with matrix flow (3) is summarized in Algorithm 1 where ALGO1 and ALGO2 could be arbitrary sub-routines producing the solution of (2) in step 0 and the following steps.

Algorithm 1 Framework to solve (2)

Inputs: Initial matrix $A^{(0)}$, vectors $\{g^{(t)}, r^{(t)}, t = 1, 2, \dots\}$.

Procedure:

- 1: Initialize: $x^{(0)} \leftarrow \text{ALGO1}(A^{(0)}, r^{(0)})$;
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: $x^{(t)} \leftarrow \text{ALGO2}(A^{(t-1)}, g^{(t)}, r^{(t)}; x^{(t-1)})$;
- 4: $A^{(t)} \leftarrow A^{(t-1)} + g^{(t)}(g^{(t)})^T$.
- 5: **end for**

Output: $\{x^{(t)} : t = 0, 1, \dots\}$.

In this article we will focus on the online part and derive an efficient algorithm for ALGO2. The complexity of ALGO1 will be increasingly less important as t increases. ALGO1 can be simply chosen as any state-of-the-art algorithm such as interior-point method. For certain problem, such as Online Newton Step [20], $A^{(0)}$ is a scaled identity matrix and hence $x^{(0)} = \frac{1}{n}\mathbf{1}$.

2.1 KKT Condition Within A Step

We start exploring the problem by considering fixed t . The aim is to minimize $\frac{1}{2}x^T Ax - r^T x$ over \mathbb{R}^n subject to $\mathbf{1}^T x = 1$ and $x \geq 0$, where A and r are abbreviation of $A^{(t)}$ and $r^{(t)}$. By Strong duality, it is equivalent to minimize the Lagrangian form

$$L(x; \mu_0, \mu) = \frac{1}{2}x^T Ax - r^T x + \mu_0(1 - \mathbf{1}^T x) - \mu^T x \quad (6)$$

where μ_0 and μ are Lagrangian multipliers with constraint $\mu_i \geq 0$ for $i = 1, \dots, n$. Denote S_x by the support of vector x . To be concise, the subscript x is suppressed in the following context. KKT condition together with Slater's condition implies that (x, μ_0, μ) is the solution of (6) if and only if

$$Ax - \mu_0 \mathbf{1} - \mu - r = 0; \quad (7)$$

$$\mathbf{1}^T x = 1; \quad (8)$$

$$\mu_i x_i = 0, \mu_i \geq 0, x_i \geq 0, \forall i = 1, \dots, n. \quad (9)$$

Here (9) is referred to as *complementary slackness* condition. The definition of $S = \text{supp}(x)$ entails that $x_{S^c} = 0$, and (9) further implies that $\mu_S = 0$. Then the condition (7) can be reformulated as

$$\begin{pmatrix} A_{SS} & A_{SS^c} \\ A_{S^cS} & A_{S^cS^c} \end{pmatrix} \begin{pmatrix} x_S \\ 0 \end{pmatrix} = \mu_0 \begin{pmatrix} \mathbf{1}_S \\ \mathbf{1}_{S^c} \end{pmatrix} + \begin{pmatrix} 0 \\ \mu_{S^c} \end{pmatrix} + \begin{pmatrix} r_S \\ r_{S^c} \end{pmatrix}$$

By separating S and S^c , we have the following equations for x_S and μ_{S^c} .

$$x_S = \mu_0 A_{SS}^{-1} \mathbf{1}_S + A_{SS}^{-1} r_S; \quad (10)$$

$$\mu_{S^c} = A_{S^cS} x_S - \mu_0 \mathbf{1}_{S^c} - r_{S^c} = -\mu_0 (\mathbf{1}_{S^c} - A_{S^cS} A_{SS}^{-1} \mathbf{1}_S) - (r_{S^c} - A_{S^cS} A_{SS}^{-1} r_S). \quad (11)$$

The other parameter μ_0 can be solved from (8) and (10). In fact,

$$1 = \mathbf{1}^T x = \mathbf{1}_S^T x_S = \mu_0 \mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S + \mathbf{1}_S^T A_{SS}^{-1} r_S$$

which implies that

$$\mu_0 = \frac{1 - \mathbf{1}_S^T A_{SS}^{-1} r_S}{\mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S}. \quad (12)$$

In summary, the quadruple $(S, x_S, \mu_{S^c}, \mu_0)$ which solves (10)-(12) produces the unique solution of (6). Moreover, given the correct support S , we can uniquely solve the other three parameters. Thus, determining S is the key part in this problem.

2.2 Homotopic Online Projection (HOP) Algorithm

2.2.1 Construction of Homotopy Continuation

Based on the above argument, the problem is reduced to updating support S with A replaced by $A + gg^T$, and r replaced by $r + \ell$, where g, ℓ are shorthand notations of $g^{(t)}$ and $r^{(t)} - r^{(t-1)}$. Heuristically, S will not be significantly disturbed when g, ℓ are small perturbations. However, in real problems, there is usually no such constraint on g . Instead, we can consider a homotopy from (A, r) to $(A + gg^T, r + \ell)$. The most natural one is $(A + \lambda gg^T, r + \underline{\lambda} \ell)$ with $\lambda, \underline{\lambda} \in [0, 1]$. In other words, if we denote $x(\lambda, \underline{\lambda})$ be the solution of (6) with (A, r) replaced by $(A + \lambda gg^T, r + \underline{\lambda} \ell)$, then $x(0, 0)$ is the solution in the last step and $x(1, 1)$ is the solution after the update. The idea of homotopy continuation method is to calculate $x(\lambda, \underline{\lambda})$ over a path linking $(0, 0)$ to $(1, 1)$. Theoretically, any path suffices and the goal is to find a path which leads to simple computation. In this article we will consider the Manhattan path from $(0, 0)$ to $(1, 1)$, namely the union of three segments: $\{(z, 0) : z \in [0, 1]\}$ and $\{(1, z) : z \in [0, 1]\}$. In other words we first minimize

$$H^{(1)}(\lambda) \triangleq \frac{1}{2} x^T (A + \lambda gg^T) x - r^T x$$

for each $\lambda \in [0, 1]$ and then minimize

$$H^{(2)}(\underline{\lambda}) \triangleq \frac{1}{2} x^T (A + gg^T) x - (r + \underline{\lambda} \ell)^T x$$

for each $\underline{\lambda} \in [0, 1]$.

Although the problem is augmented, the update is efficient since the support S is shown to be a piecewise constant set on the path. The explicit formulas, namely (10) - (12), can be used to compute (x_S, μ_{S^c}, μ_0) directly when S is fixed. In fact, the triple (x_S, μ_{S^c}, μ_0) is a simple function of $(\lambda, \underline{\lambda})$ as shown in the following theorem.

Theorem 1

1. For a given $\underline{\lambda}$, there exist vectors $u_1, u_2 \in \mathbb{R}^{n+1}$ and scalars $D_1, D_2 \in \mathbb{R}$, which only depend on S , such that

$$\begin{pmatrix} x_S(\lambda) \\ -\mu_{S^c}(\lambda) \\ \mu_0(\lambda) \end{pmatrix} = \frac{u_1 - u_2 \lambda}{D_1 - D_2 \lambda}. \quad (13)$$

2. For given λ , there exist vectors $\underline{u}_1, \underline{u}_2 \in \mathbb{R}^{n+1}$, which only depend on S , such that

$$\begin{pmatrix} x_S(\underline{\lambda}) \\ -\mu_{S^c}(\underline{\lambda}) \\ \mu_0(\underline{\lambda}) \end{pmatrix} = \underline{u}_1 - \underline{u}_2 \underline{\lambda}. \quad (14)$$

Proof

1. The proof is quite involved and we relegate it into Theorem 4 in Appendix B. The theorem also gives the exact formula of u_1, u_2, D_1, D_2 .

2. By (12), we have

$$\mu_0(\lambda) = \frac{1 - \mathbf{1}_S^T A_{SS}^{-1} (r_S + \lambda \ell_S)}{\mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S} = \mu_0(0) - \frac{\mathbf{1}_S^T A_{SS}^{-1} \ell_S}{\mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S} \lambda.$$

Then it follows from (10) that

$$x_S(\lambda) = \mu_0(\lambda) A_{SS}^{-1} \mathbf{1}_S + A_{SS}^{-1} (r_S + \lambda \ell_S) = x_S(0) - \left(\frac{\mathbf{1}_S^T A_{SS}^{-1} \ell_S}{\mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S} A_{SS}^{-1} \mathbf{1}_S - A_{SS}^{-1} \ell_S \right) \lambda.$$

Similarly, by (11), we obtain that

$$-\mu_{S^c}(\lambda) = -\mu_{S^c}(0) - \left(\frac{\mathbf{1}_S^T A_{SS}^{-1} \ell_S}{\mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S} (\mathbf{1}_{S^c} - A_{S^c S} A_{SS}^{-1} \mathbf{1}_S) - (\ell_{S^c} - A_{S^c S} A_{SS}^{-1} \ell_S) \right) \lambda.$$

■

2.2.2 Update of Support

Once $S = S(\lambda)$ is obtained for all λ , the solution path can be efficiently solved by Theorem 1. Heuristically, S is piecewise constant and the task is reduced to finding the next λ that $S(\lambda)$ changes. We consider the update of S in optimizing $H^{(1)}(\lambda)$. The update of S in optimizing $H^{(2)}(\lambda)$ can be obtained in the same way.

For a given $\lambda_0 \in [0, 1]$, if $x_S(\lambda_0) > 0$ and $\mu_{S^c}(\lambda_0) > 0$, then (13) implies that there exists $\eta > 0$, such that for any $\lambda \in (\lambda_0 - \eta, \lambda_0 + \eta)$, both $x_S(\lambda)$ and $\mu_{S^c}(\lambda)$ remain positive by setting $S(\lambda) = S(\lambda_0)$. Since (10)-(12) are sufficient and necessary, we conclude that $S(\lambda) = S(\lambda_0)$. This argument remains valid until an entry of either x_S or μ_{S^c} hits zero. Denote j by the index of this entry. In the former case, j leaves S and S is updated to $S \setminus \{j\}$. In the latter case, j enters into S and S is updated to $S \cup \{j\}$. The other three parameters are then updated correspondingly by Theorem 1. Theorem 2 formalizes the above claim. The proof is omitted since it is a direct consequence of sufficiency and necessity of KKT conditions (10)-(12).

Theorem 2 For any given $\lambda_0 \geq 0$, let λ^{new} be the next smallest λ such that one entry of either x_S or μ_{S^c} hits 0, i.e.

$$\lambda^{\text{new}} = \min_+ \left\{ \frac{u_{1i}}{u_{2i}} : i = 1, 2, \dots, n \right\},$$

where u_1 and u_2 are defined in (13) and \min_+ evaluates the minimum positive number in the set and defined to be ∞ if all elements are non-positive. Then $S(\lambda) \equiv S(\lambda_0)$ for $\lambda \in [\lambda_0, \lambda^{\text{new}}]$. Further, let

$$I_1 = \{i \in S : u_{1i} = u_{2i} \lambda^{\text{new}}\}, \quad I_2 = \{i \in S^c : u_{1i} = u_{2i} \lambda^{\text{new}}\},$$

then $S(\lambda^{\text{new}})$ is updated by

$$S(\lambda^{\text{new}}) = (S(\lambda_0) \setminus I_1) \cup I_2.$$

Remark 1 According to our experience, $I_1 \cup I_2$ at most contains one element. In other words, S is updated by one element at each time.

In summary, the algorithm starts from $\lambda = 0$ and searches for the next smallest λ such that one entry of x_S or μ_{S^c} hits zero, then updates λ as well as the quadruple $(S, x_S, \mu_{S^c}, \mu_0)$. The procedure is repeated until λ crosses 1. In other words, there exist a sequence $0 = \lambda_0 < \lambda_1 < \dots < \lambda_k = 1$, which we call *turning points*, such that $x(\lambda)$ has the same support between any two consecutive turning points and the value can be calculated by Theorem 1. A counterpart of Theorem 2 can be established for λ . The whole task amounts to finding all turning points and we call this procedure *homotopic online projection (HOP)* algorithm. The complexity of the HOP algorithm is determined by both the number of turning points and the complexity of the update between two consecutive turning points. To be more clear, we state the main steps in Algorithm 2 for optimizing $H^{(1)}(\lambda)$, i.e. with $\lambda = 0$. For compact notation, we define v as a $n \times 1$ vector with $v_S = x_S$ and $v_{S^c} = \mu_{S^c}$. As a convention, the minimum of an empty set is set to be infinity (line 3 of Algorithm 2).

Algorithm 2 Main steps of the HOP algorithm in optimizing $H^{(1)}(\lambda)$

Inputs: parameters A, y, r, g ; initial optimum x (corresponding to A)

Procedure:

```
1: Initialize  $\lambda \leftarrow 0, S \leftarrow \text{supp}(x)$ ;  
2: while  $\lambda < 1$  do  
3:    $\lambda = \min\{\lambda_1 > \lambda : v_i(\lambda_1) = 0 \text{ for some } i\}$ ;  
4:    $I_1 \leftarrow \{i \in S : v_i(\lambda) = 0\}$ ;  
5:    $I_2 \leftarrow \{i \in S^c : v_i(\lambda) = 0\}$ ;  
6:   if  $\lambda \leq 1$  then  
7:      $S \leftarrow (S \setminus I_1) \cup I_2$ ;  
8:   else  
9:      $\lambda \leftarrow 1$ ;  
10:  end if  
11:   $(x_S, \mu_{S^c}, \mu_0) \leftarrow (x_S(\lambda), \mu_{S^c}(\lambda), \mu_0(\lambda))$  via (10)-(12).  
12: end while
```

Output: $(S, x_S, \mu_{S^c}, \mu_0)$.

2.3 Implementation and Complexity Analysis

Algorithm 2 presents the main idea without the implementation details. Although we can implement Algorithm 2 by directly computing quantities, e.g. u_1, u_2 , in every step to find the next turning point as in line 3 and also directly computing the iterates via (10)-(12) as in line 11, it is fairly inefficient since many quantities appear in several computation steps and we can store them to save the computation. A careful derivation in Appendices B and C shows that the computation complexity is indeed low. Theorem 3 summarizes the complexity for optimizing $H^{(1)}(\lambda), H^{(2)}(\underline{\lambda})$ separately. As a convention, we assume the scalar-scalar multiplication takes a unit time and ignore the addition for simplicity. Since the real implementation is involved, we state it as well as the proof of theorem 3 in Appendices B and C for two cases separately.

Theorem 3 *In step t , denote by k_A, k_r the number of turning points in optimizing $H^{(1)}(\lambda)$ and $H^{(2)}(\underline{\lambda})$. Further let s be the maximum support size over the path of $(\lambda, \underline{\lambda})$ and s_* by the size of union of all supports from step 1 to step t . Let C_{jt} be the computation cost of the HOP algorithm in optimizing $H^{(j)}$, then*

1. $C_{1t} = ns_* + ns(3k_A + 1) + n(12k_A + 2) + O(k_A)$;
2. $C_{2t} = ns(2k_r + 1) + n(6k_r + 1) + O(k_r)$.

It is clear that the algorithm adapts to the sparsity when optimizing both $H^{(1)}(\lambda)$ and $H^{(2)}(\underline{\lambda})$. According to our experience, the solution is usually extremely sparse compared to the dimension. For example, in section 3 we show that the average support size of NASDAQ data from Jan. 3, 2005 to May. 13, 2016 is 160 among 1544 stocks and that of NYSE data in the same time range is 89 among 1101 stocks. More generally, the solution sparsity of quadratic programming problem is observed in many other applications and has been theoretically justified when A has certain structures [6, 10].

2.4 Number of Turning Points

Let S_t be the support of the optimum and k_t be the number of total turning points, then we can derive a generic bound that

$$k_t \geq |S_t \setminus S_{t-1}| + |S_{t-1} \setminus S_t| \quad (15)$$

provided that only one element is added to or removed from the support at each update; see Remark 1. This is because it requires at least $|S_{t-1} \setminus S_t|$ steps to pop out the elements in $S_{t-1} \setminus S_t$ and $|S_t \setminus S_{t-1}|$ steps to push in the elements in $S_t \setminus S_{t-1}$ to translate S_{t-1} into S_t .

On the other hand, suppose that no other coordinates than those in $S_t \cup S_{t-1}$ enter into the support in the path, then

$$k_t = |S_t \setminus S_{t-1}| + |S_{t-1} \setminus S_t|. \quad (16)$$

Heuristically, the equation (16) should hold since if a coordinate, not in $S_t \cup S_{t-1}$, entered into the support on the path, it must be popped out before the end but this should be rare to happen. For both synthetic data and real data in section 3, we observed that there are at least 95% of steps with k_t satisfying (16) and over 99% of steps with $k_t \leq |S_t \setminus S_{t-1}| + |S_{t-1} \setminus S_t| + 6$, i.e. with at most 3 outside coordinates entered into the path. Thus, (16) is a highly reliable result for k_t .

As a direct consequence of (16), the HOP algorithm is efficient when the support changes slowly so that k_t is small. In addition, if the solution is sparse, then a rough bound suggests that $k_t \leq |S_t| + |S_{t-1}|$ is small. These phenomena were observed in various situations (see section 3) and (16) explains the good performance of the HOP algorithm.

3 Experiments

In this section, we compared the performance of HOP with SPG and PQN on both synthetic data and real data examples. We implemented HOP in MATLAB³ and implemented SPG and PQN by using existing code⁴. We also tested the interior point method using the function *quadprog* but did not report the result in our figures and tables since it is far too slow, though we still mentioned the result in the following subsections. It is not surprising since interior-point method has a $O(n^3T)$ complexity in each step. To evaluate the performance, we displayed the running time as a measure of efficiency.

3.1 Synthetic Data

We consider the following problem with synthetic data:

$$\min_{x \in \Delta_n} \frac{1}{2}(x - y)^T A^{(t)}(x - y). \quad (17)$$

Without the superscript t , this problem is called *standard quadratic programming problem* and has attracted the attention in various fields, e.g. [6, 26, 7]. It is of particular interest to study the case where $A^{(t)}$ is a random matrix generated from some distribution. For instance, [10] consider a Wigner matrix A with $\{A_{ij} : 1 \leq i \leq j \leq n\}$ being i.i.d. random variables and $A_{ji} = A_{ij}$. In this article, we consider another important class of matrices in random matrix theory — covariance matrix of rectangular matrices with i.i.d. entries, i.e. $A^{(t)} = \sum_{s=1}^t g^{(s)}(g^{(s)})^T$, where $g^{(s)} \in \mathbb{R}^n$ has i.i.d. Gaussian entries; see [3] for more discussion. In this case the matrix flow $\{A^{(t)} : t = 1, 2, \dots\}$ satisfies (3). To avoid singularity, we set $A^{(0)} = \epsilon I$ where $\epsilon = 10^{-4}$ is a small positive number. Then it is easy to see that $A^{(t)}$ is non-singular for all t with probability 1 so that the solution $x^{(t)}$ is unique. The vector y governs the sparsity of the solution. To see this, consider the isotropic case where $A = I$, the solution of (17) is the projection of y onto the simplex. If y is a zero vector, the optimum is a dense vector with all entries $\frac{1}{n}$. In contrast, if y has large entries, the simplex constraint will pull the optimum towards that direction and force the other entries to be zero, in which case the solution is sparse. The same phenomenon was observed in anisotropic case as shown below.

Our goal is to explore the scalability, in terms of the dimension, and the adaptivity to solution sparsity of the algorithms. For the former we considered three dimensions: $\{100, 1000, 3000\}$ and for the latter we set $y = cy_0$ with y_0 generated from $N(0, I_{n \times n})$ and $c \in \{0.01, 0.1\}$. For each case, we set the total number of steps as 5000 and treated every 250 steps as an epoch (20 epochs in total). In each epoch, we took the HOP algorithm as a benchmark and reported the ratio of the running time of other algorithms and that of the HOP algorithm on a machine with 2 GHz Intel Core i7 processor. For instance, a ratio 2 for SPG means that HOP is twice as fast as SPG in the given epoch. To address the relationship between solution sparsity and computation efficiency, we reported the average size of supports as well as its standard deviation and maximum in Table 1. In addition, we reported the ratio of the overall running time of SPG/PQN and that of HOP as a summarized measure for the computation gain. As expected, a larger c gives sparser solutions along the path and HOP significantly outperforms SPG when the solution is sparse ($c = 0.1$) especially for large-scale problems ($n = 1000, 3000$), in which HOP is over 4.5 times faster than SPG. When the solution is not sparse ($c = 0.01$), HOP is similar to SPG in small-scale problem ($n = 100$) and increasingly

³Code available at <https://github.com/Elric2718/HOP>.

⁴<https://www.cs.ubc.ca/~schmidt/Software/minConf.html>

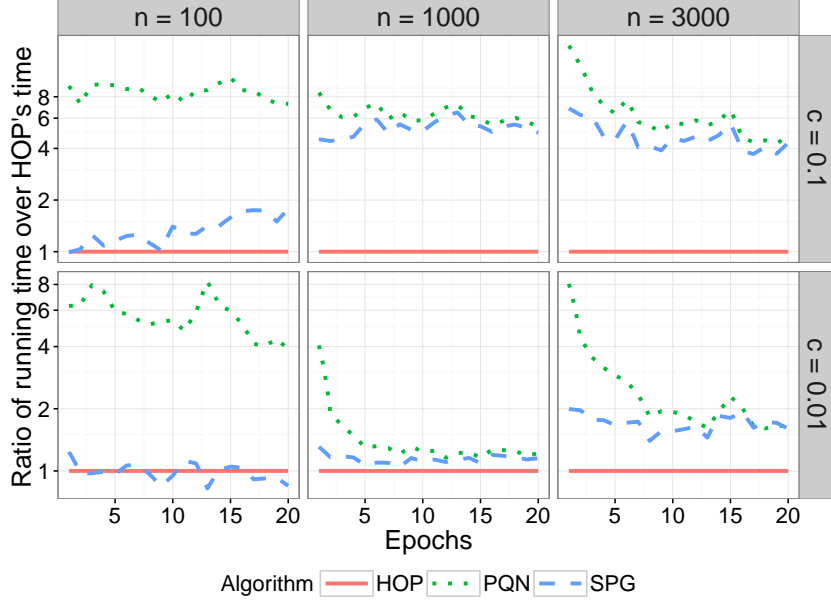


Figure 1: Ratio of the average execution time of SPG/PQN and that of HOP evaluated every 250 steps, on synthetic datasets (on a machine with 2 GHz Intel Core i7 processor). Each column corresponds to a dimension n and each row corresponds to a factor c .

more efficient when the size of the problem grows. In all cases, HOP is much more efficient than PQN. In summary, HOP is more scalable and more adaptive to the solution sparsity than SPG and PQN.

scenarios		sparsity		speed ratio	
n	c	mean (std.)	max	SPG	PQN
100	0.01	81.3 (3.5)	88	0.97	5.72
1000	0.01	158.3 (30.6)	190	1.13	1.31
3000	0.01	146.4 (34.9)	225	1.68	2.24
100	0.1	19.5 (2.3)	26	1.34	8.47
1000	0.1	22.0 (5.3)	32	5.26	6.20
3000	0.1	18.3 (4.2)	27	4.51	5.88

Table 1: Solution Sparsity and overall computation gain of HOP over SPG/PQN on synthetic datasets. The first two columns correspond to the dimension and the factor c ; the third column gives the mean of support size with its standard deviation (in the parentheses); the fourth column gives the maximum support size along the path; the last two columns show the ratio of overall running time between SPG/PQN and HOP.

Finally, we tested our conjecture in section 2.4 on the number of turning points. As explained there, a benchmark for k_t is $|S_t \setminus S_{t-1}| + |S_{t-1} \setminus S_t|$. We refer to $e_t = (k_t - |S_t \setminus S_{t-1}| + |S_{t-1} \setminus S_t|)/2$ as the number of *excess turning points*; see section 2.4 for details. For each synthetic dataset, we reported the proportion of zero e_t in Table 2. It is clear that $e_t = 0$ for over 98.5% of the steps and over 99.9% of the steps have $e_t \leq 6$ in all scenarios. This validates our heuristics on the number of turning points.

3.2 Real Data

In this part, we considered the application of the HOP algorithm on the Markowitz portfolio selection problem with sequential data, formulated as (4). The most natural estimators $\hat{\Sigma}^{(t)}$ and $\hat{\mu}^{(t)}$ are the

scenarios		proportion of zeros	quantiles	
n	c		99%	99.9%
100	0.01	99.8%	0	1
1000	0.01	98.9%	1	4
3000	0.01	98.7%	1	6
100	0.1	99.9%	0	0
1000	0.1	99.8%	0	1
3000	0.1	99.8%	0	1

Table 2: Distribution of e_t , the number of excess turning points, on synthetic datasets. The first two columns correspond to the dimension and the factor c ; the third column gives the proportion of zero e_t ; the last two columns give the 99% and 99.9% quantiles of e_t .

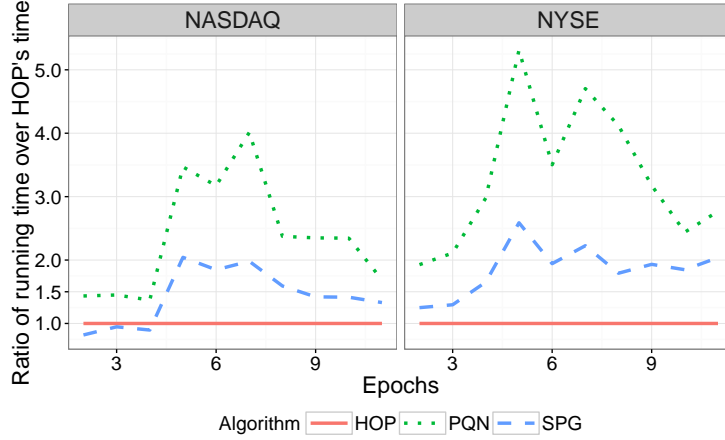


Figure 2: Ratio of the average execution time of SPG/PQN and HOP, evaluated every 252 steps, on NYSE and NASDAQ datasets (on a machine with 2 GHz Intel Core i7 processor).

sample covariance matrix and the sample mean, i.e.

$$\hat{\mu}^{(t)} = \frac{1}{t} \sum_{s=1}^t w^{(s)}, \hat{\Sigma}^{(t)} = \frac{1}{t} \sum_{s=1}^t (w^{(s)} - \hat{\mu}^{(s)})(w^{(s)} - \hat{\mu}^{(s)})^T,$$

where $w^{(t)}$ is the vector of daily gains, measured by the log price ratio $\log(p_t/p_{t-1})$, of all assets of interest. Via some algebra, it can be shown that the flow $\{t\Sigma^{(t)} : t = 1, 2, \dots\}$ satisfies (3) with

$$\begin{aligned} t\hat{\Sigma}^{(t)} &= (t-1)\hat{\Sigma}^{(t-1)} + \frac{t-1}{t} \left(w^{(t)} - \hat{\mu}^{(t-1)} \right) \left(w^{(t)} - \hat{\mu}^{(t-1)} \right)^T, \\ t\hat{\mu}^{(t)} &= (t-1)\hat{\mu}^{(t-1)} + w^{(t)}. \end{aligned}$$

Thus, the problem (4) is equivalent to the problem (2) + (3) with

$$A^{(t)} = t\hat{\Sigma}^{(t)}, r^{(t)} = t\hat{\mu}^{(t)}, g^{(t)} = \sqrt{\frac{t-1}{t}} \left(w^{(t)} - \hat{\mu}^{(t-1)} \right).$$

Our goal is to calculate the optimum of (4) for each step t . We should emphasize that the solution in this way is optimal from hindsight, which is different from the notions in online learning community. Nonetheless, it is an interesting and important problem in the context of back testing and risk management since the result can reveal the hidden structure of the assets; see [8, 15, 16] for more details.

Now we considered two datasets from NYSE and NASDAQ⁵, with time range from Jan. 3, 2005 to May. 13, 2016. The NYSE data contains 1544 stocks and NASDAQ data contains 1101 stocks.

⁵Data available at <https://github.com/Elric2718/HOP>.

This is in contrast to the classical studies where at most hundreds of stocks, e.g. S&P500 are incorporated. However, we should emphasize that for some financial institutions like hedge funds, the number of base assets is huge and the computation efficiency becomes important when the trading frequency is high. Here we consider a large number of stocks to show the potential of the HOP algorithm in optimizing a large basket of assets. Another minor issue is shorting. Although allowed in U.S. market, the short sale is much harder to operate than the long sale and is usually under more strict regulation. In practice, it is reasonable to assume that all weights $x_i \geq -\underline{x}$ for some positive \underline{x} . We can transform x_i into \tilde{x}_i by $\tilde{x}_i = \frac{x_i + \underline{x}}{1 + n\underline{x}}$ and transform the optimization problem (4) accordingly. Then the new vector \tilde{x} satisfies the simplex constraint. However, the main concern of this article is the computation efficiency and hence we consider the case where shorting is disallowed. For convenience, we set $\eta = 0$ in (4).

Similar to the previous subsection, we reported the ratio of execution time in Figure 2 every 252 steps (there are 252 trading days in most years) and reported other results in Table 3 and Table 4. It is clearly seen from Figure 2 that the HOP algorithm outperforms SPG and PQN on both datasets and the computation gain of the HOP algorithm on NYSE dataset is more significant than that on NASDAQ dataset since the solution is more sparse, as shown in Table 3. In the same table, we reported the overall computation gain as well as the computation gain after 4 epochs. In both settings, HOP is faster than SPG and PQN and the gain is more significant in the later stages since HOP becomes more stable after a certain number of steps so that the number of turning points decreases significantly. We also tried the interior-point method on both datasets. It is 67 times slower than HOP on NYSE dataset and 31 times slower than HOP on NASDAQ dataset. We did not report it in the figure since it is too slow. Finally, the heuristics on the number of turning points is also validated on these datasets as shown in Table 4.

Dataset	sparsity		speed ratio (s.r.)		s.r. (after 4 epochs)	
	mean (std.)	max (min)	SPG	PQN	SPG	PQN
NYSE	49.8 (22.7)	108 (30)	1.81	3.21	2.04	3.70
NSADAQ	148.9 (17.8)	192 (127)	1.37	2.26	1.67	2.78

Table 3: Solution Sparsity and the computation gain of HOP on NYSE and NASDAQ datasets. The former includes the mean, standard deviation, maximum and minimum support size; the latter includes the ratio of overall execution time and the ratio of execution time after 4 epochs between SPG/PQN and HOP.

Dataset	proportion of zeros	quantiles	
		99%	99.9%
NYSE	97.9%	1	10
NSADAQ	96.2%	3	22

Table 4: distribution of e_t , the number of excess turning points, on NYSE and NASDAQ datasets.

4 Discussion

In this article, we proposed an efficient algorithm (HOP) to solve a sequential quadratic programming problems (2) with rank-one update (3) based on a homotopy continuation that links the consecutive objectives. By a careful derivation, we calculated the exact complexity in each step up to a universal constant (Theorem 3) and showed that the HOP algorithm has a good performance when the support of the solution changes slowly with time or is sparse as in many applications. The efficiency of the HOP algorithm was validated on both synthetic data and real data.

As mentioned in section 1, another example that involves the generic problem (2) is the Online Newton Step, which proved to be effective in the context of online learning [20]. However, the projection step (5) is the bottleneck which dominates the other steps when using projected gradient descent method, as suggested by [1], or the interior point method, as suggested by [20], and the computation burden limits the exploration of ONS. However, we believe that the HOP algorithm can

solve (2) much more efficiently than existing methods and ONS can attract more attention once the computation bottleneck is broken.

Another interesting direction is to apply the HOP algorithm on offline problem, i.e. minimizing (2) with a single A and r . Although no online movement exists, we can artificially define $r^{(0)}$ as a zero vector and $A^{(0)}$ as a matrix with n -th row and n -th column equal to 0 such that $A = A^{(0)} + gg^T$ for some $g \in \mathbb{R}^n$. For instance, it is not hard to show that the particular choice $g = \frac{Ae_n}{A_{nn}}$ gives the desired result, where e_n is the n -th basis vector in \mathbb{R}^n with n -th entry equal to 1 and all others equal to 0. Then $x^{(0)} = e_n$ is the optimum corresponding to $A^{(0)}$ and $r^{(0)}$. We can start from $x^{(0)}$ and apply the HOP algorithm for one step to solve the problem.

Finally, the simplex constraint is not essential in the HOP algorithm and it can be generalized to deal with more linear and quadratic constraints via a more involved derivation. It would be an important future direction to write a generic solver for sequential quadratic programming problems with general linear and quadratic constraints.

References

- [1] A. Agarwal, E. Hazan, S. Kale, and R. E. Schapire. Algorithms for portfolio management based on the newton method. In *Proceedings of the 23rd international conference on Machine learning*, pages 9–16. ACM, 2006.
- [2] E. L. Allgower and K. Georg. *Numerical continuation methods: an introduction*, volume 13. Springer Science & Business Media, 2012.
- [3] Z. Bai and J. W. Silverstein. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.
- [4] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [5] E. G. Birgin, J. M. Martínez, and M. Raydan. Nonmonotone spectral projected gradient methods on convex sets. *SIAM Journal on Optimization*, 10(4):1196–1211, 2000.
- [6] I. M. Bomze. On standard quadratic optimization problems. *Journal of Global Optimization*, 13(4):369–387, 1998.
- [7] I. M. Bomze, M. Locatelli, and F. Tardella. New and old bounds for standard quadratic optimization: dominance, equivalence and incomparability. *Mathematical Programming*, 115(1):31–64, 2008.
- [8] J. Brodie, I. Daubechies, C. De Mol, D. Giannone, and I. Loris. Sparse and stable markowitz portfolios. *Proceedings of the National Academy of Sciences*, 106(30):12267–12272, 2009.
- [9] N. Cesa-Bianchi, Y. Freund, D. Haussler, D. P. Helmbold, R. E. Schapire, and M. K. Warmuth. How to use expert advice. *Journal of the ACM (JACM)*, 44(3):427–485, 1997.
- [10] X. Chen, J. Peng, and S. Zhang. Sparse solutions to random standard quadratic optimization problems. *Mathematical Programming*, 141(1-2):273–293, 2013.
- [11] S.-N. Chow, J. Mallet-Paret, and J. A. Yorke. A homotopy method for locating all zeros of a system of polynomials. In *Functional differential equations and approximation of fixed points*, pages 77–88. Springer, 1979.
- [12] T. M. Cover. Universal portfolios. *Mathematical finance*, 1(1):1–29, 1991.
- [13] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the l_1 -ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pages 272–279. ACM, 2008.
- [14] B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [15] J. Fan, J. Zhang, and K. Yu. Asset allocation and risk assessment with gross exposure constraints for vast portfolios. *Available at SSRN 1307423*, 2008.
- [16] J. Fan, J. Zhang, and K. Yu. Vast portfolio selection with gross-exposure constraints. *Journal of the American Statistical Association*, 107(498):592–606, 2012.

- [17] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- [18] P. Garrigues and L. E. Ghaoui. An homotopy algorithm for the lasso with online observations. In *Advances in neural information processing systems*, pages 489–496, 2009.
- [19] E. Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- [20] E. Hazan, A. Kalai, S. Kale, and A. Agarwal. Logarithmic regret algorithms for online convex optimization. In *International Conference on Computational Learning Theory*, pages 499–513. Springer, 2006.
- [21] T. Ibaraki and N. Katoh. *Resource allocation problems: algorithmic approaches*. MIT press, 1988.
- [22] T.-Y. Li. On chow, mallet-paret and yorke homotopy for solving systems of polynomials. *Bull. Inst. Math. Acad. Sinica*, 11(3):433–437, 1983.
- [23] M. P. F. M. Schmidt, E. van den Berg and K. Murphy. Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm. In *Proceedings of The Twelfth International Conference on Artificial Intelligence and Statistics (AISTATS) 2009*, volume 5, pages 456–463, 2009.
- [24] H. Markowitz. Portfolio selection. *The journal of finance*, 7(1):77–91, 1952.
- [25] H. Markowitz. The optimization of a quadratic function subject to linear constraints. *Naval research logistics Quarterly*, 3(1-2):111–133, 1956.
- [26] A. Scozzari and F. Tardella. A clique algorithm for standard quadratic programming. *Discrete Applied Mathematics*, 156(13):2439–2448, 2008.
- [27] V. G. Vovk. A game of prediction with expert advice. In *Proceedings of the eighth annual conference on Computational learning theory*, pages 51–60. ACM, 1995.
- [28] P. Wolfe. The simplex method for quadratic programming. *Econometrica: Journal of the Econometric Society*, pages 382–398, 1959.
- [29] S. J. Wright. *Primal-Dual Interior-Point Methods*, volume 54. SIAM, 1997.

A Roadmap of Appendices

The general idea of HOP algorithm has been presented in section 2. However, to implement it efficiently, we need much more efforts to explore the structure of the solution path and find common quantities which are used by multiple sub-routines. To make the derivation well-organized, we start from considering the case where only one of A and r is time-varying while the other parameter is fixed. The case where A is time-varying and the case where r is time-varying are considered separately in Appendix B and C. Then in Appendix D, we combine two components and state the implementation for the general case.

In each following appendix, we will first define a list of case-specific intermediate variables, which are the key ingredients to improve efficiency. Then we describe the whole procedure followed by details of each sub-routine. Finally, we give a complexity analysis at the end of each appendix.

B Implementation of HOP Algorithm With Time-Varying A and Fixed r

B.1 Intermediate Variables

Although (10)-(12) completely define the solution, they involve messy terms. To simplify the notations, we define a list of intermediate variables. These variables play important roles in the implementation since they capture the quantities repeatedly appeared and unnecessary computation can be avoided by storing their values in memory.

The intermediate variables are defined as follows. First, let M be a $n \times n$ matrix such that

$$M_{SS} = A_{SS}^{-1}, \quad M_{S^c S} = -A_{S^c S} A_{SS}^{-1}, \quad M_{\cdot, S^c} = 0. \quad (18)$$

For large-scale problem where n is prohibitively large, we can only store a $n \times |S|$ matrix by removing the zero entries of M . This saves storage cost significantly. Then we define two vectors $\eta, \tilde{\eta} \in \mathbb{R}^n$ such that

$$\eta_S = M_{SS}g_S, \quad \eta_{S^c} = g_{S^c} + M_{S^cS}g_S, \quad \tilde{\eta}_S = M_{SS}\mathbf{1}_S, \quad \tilde{\eta}_{S^c} = \mathbf{1}_{S^c} + M_{S^cS}\mathbf{1}_S. \quad (19)$$

Last we define four scalars.

$$D = \mathbf{1}_S^T A_{SS}^{-1} \mathbf{1}_S, \quad D_g = \mathbf{1}_S^T A_{SS}^{-1} g_S, \quad D_{gg} = g_S^T A_{SS}^{-1} g_S, \quad D_{gr} = -\eta_S^T r_S. \quad (20)$$

Note that all variables are functions of λ if A is replaced by $A + \lambda g g^T$ and we denote them by $\bullet(\lambda)$. For example,

$$D(\lambda) = \mathbf{1}_S^T (A_{SS} + \lambda g_S g_S^T)^{-1} \mathbf{1}_S,$$

and others can be defined in a similar fashion. The following lemma formulates these functions.

Lemma 1 Let $\alpha(\lambda) = \frac{\lambda}{1+\lambda D_{gg}}$. Before any entry of (x_S, μ_{S^c}) hits 0, it holds that

- $M_{\cdot,S}(\lambda) = M_{\cdot,S} - \alpha(\lambda)\eta\eta_S^T$;
- $\eta(\lambda) = \frac{\eta}{1+\lambda D_{gg}}$;
- $\tilde{\eta}(\lambda) = \tilde{\eta} - \alpha(\lambda)D_g\eta$;
- $D(\lambda) = D - \alpha(\lambda)D_g^2$;
- $(D_g(\lambda), D_{gg}(\lambda), D_{gr}(\lambda)) = \frac{1}{1+\lambda D_{gg}}(D_g, D_{gg}, D_{gr})$.

Proof By Sherman-Morrison-Woodbury formula,

$$M_{SS}(\lambda) = (A_{SS} + \lambda g_S g_S^T)^{-1} = A_{SS}^{-1} - \lambda \frac{A_{SS}^{-1} g_S g_S^T A_{SS}^{-1}}{1 + \lambda g_S^T A_{SS}^{-1} g_S} = M_{SS} - \alpha(\lambda)\eta_S \eta_S^T.$$

This implies that

$$\begin{aligned} M_{S^cS}(\lambda) &= -(A_{S^cS} + \lambda g_{S^c} g_S^T)(A_{SS}^{-1} - \alpha(\lambda)\eta_S \eta_S^T) \\ &= M_{S^cS} - \lambda g_{S^c} g_S^T A_{SS}^{-1} + \lambda \alpha(\lambda) g_{S^c} g_S^T \eta_S \eta_S^T + \alpha(\lambda) A_{S^cS} \eta_S \eta_S^T \\ &= M_{S^cS} - \lambda g_{S^c} \eta_S^T + \lambda \alpha(\lambda) D_{gg} g_{S^c} \eta_S^T + \alpha(\lambda) A_{S^cS} \eta_S \eta_S^T & [Use \ D_{gg} = g_S^T \eta_S] \\ &= M_{S^cS} - (\lambda - \lambda \alpha(\lambda) D_{gg}) g_{S^c} \eta_S^T + \alpha(\lambda) A_{S^cS} \eta_S \eta_S^T \\ &= M_{S^cS} - \alpha(\lambda) g_{S^c} \eta_S^T + \alpha(\lambda) A_{S^cS} \eta_S \eta_S^T & [Use \ \lambda - \lambda \alpha(\lambda) D_{gg} = \alpha(\lambda)] \\ &= M_{S^cS} - \alpha(\lambda) (g_{S^c} - A_{S^cS} A_{SS}^{-1} g_S) \eta_S^T \\ &= M_{S^cS} - \alpha(\lambda) \eta_{S^c} \eta_S^T. \end{aligned}$$

Putting pieces together, we obtain that

$$M_{\cdot,S}(\lambda) = M_{\cdot,S} - \alpha(\lambda)\eta\eta_S^T.$$

Based on $M_{\cdot,S}(\lambda)$, it is straightforward to derive other variables. For $\eta(\lambda)$,

$$\begin{aligned} \eta_S(\lambda) &= M_{SS}(\lambda)g_S = \eta_S - \alpha(\lambda)\eta_S \eta_S^T g_S \\ &= (1 - \alpha(\lambda)D_{gg})\eta_S = \frac{\eta_S}{1 + \lambda D_{gg}}; \\ \eta_{S^c}(\lambda) &= g_{S^c} + M_{S^cS}(\lambda)g_S = \eta_{S^c} - \alpha(\lambda)\eta_{S^c} \eta_S^T g_S \\ &= (1 - \alpha(\lambda)D_{gg})\eta_{S^c} = \frac{\eta_{S^c}}{1 + \lambda D_{gg}}. \end{aligned}$$

Thus,

$$\eta(\lambda) = \frac{\eta}{1 + \lambda D_{gg}}.$$

Similarly,

$$\tilde{\eta}_S(\lambda) = M_{SS}(\lambda)\mathbf{1}_S = \tilde{\eta}_S - \alpha(\lambda)\eta_S \eta_S^T \mathbf{1}_S$$

$$\begin{aligned}
&= \tilde{\eta}_S - \alpha(\lambda) D_g \eta_S; \\
\tilde{\eta}_{S^c}(\lambda) &= \mathbf{1}_{S^c} + M_{S^c S}(\lambda) \mathbf{1}_S = \mathbf{1}_{S^c} - \alpha(\lambda) \eta_{S^c} \eta_S^T \mathbf{1}_S \\
&= \tilde{\eta}_{S^c} - \alpha(\lambda) D_g \eta_{S^c},
\end{aligned}$$

and hence

$$\tilde{\eta}(\lambda) = \tilde{\eta} - \alpha(\lambda) D_g \eta.$$

The last four scalars are even easier to handle. In fact, $D(\lambda)$ can be derived directly by

$$D(\lambda) = \mathbf{1}_S^T (M_{SS} - \alpha(\lambda) \eta_S \eta_S^T) \mathbf{1}_S = D - \alpha(\lambda) D_g^2$$

By reformulating the other three variables, the last statement can be proved,

$$\begin{aligned}
(D_g(\lambda), D_{gg}(\lambda), D_{gr}(\lambda)) &= (\mathbf{1}_S^T \eta_S(\lambda), g_S^T \eta_S(\lambda), -r_S^T \eta_S(\lambda)) \\
&= \frac{1}{1 + \lambda D_{gg}} (\mathbf{1}_S^T \eta_S, g_S^T \eta_S, -r_S^T \eta_S) \\
&= \frac{1}{1 + \lambda D_{gg}} (D_g, D_{gg}, D_{gr}).
\end{aligned}$$

■

B.2 Implementation

Lemma 1 implies that given the function values of the intermediate variables at $\lambda = 0$, the function values at a neighborhood of 0 can be calculated directly. Within time t , all intermediate variables will be updated correspondingly when the support changes. It has been shown in Lemma 1 that updating M requires $n|S|$ operations while updating other variables only requires n operations. When the problem transits from time t to time $t + 1$, the variables $(\eta, D_g, D_{gg}, D_{gr})$ needs to be recalculated since it depends on a new $g^{(t+1)}$. In contrast, $(M, \tilde{\eta}, D)$ can be updated in the same way as in time t . In summary, $(M, \tilde{\eta}, D)$ is shared by for all times while $(\eta, D_g, D_{gg}, D_{gr})$ is only used in a single time. For compact notation, we define Par_1 and Par_2 as

$$\text{Par}_1 = \{M, \tilde{\eta}, D\}, \quad \text{Par}_2 = \{\eta, D_g, D_{gg}, D_{gr}\}. \quad (21)$$

In addition, we denote v by the concatenation of x_S and $-\mu_{S^c}$, i.e.

$$v_S = x_S, \quad v_{S^c} = -\mu_{S^c}, \quad (22)$$

as a $n \times 1$ vector. It will be shown in the next subsection that $v(\lambda)$ can be expressed in a concise way.

Algorithm 3 describes the full implementation of HOP algorithm, which solves the online problem (2) with r fixed. The sub-routines involved will be discussed separately in following subsections. Roughly speaking, after initialization, we enter into the outer-loop and try to solve (3) at time t using the information from time $t - 1$. Starting from $\lambda = 0$, we search for the next λ that pushes one entry of v to zero. FIND_LAMBDA fulfills this goal and also reports the corresponding entry j . If $j \in S$ then j is removed from S and otherwise j is added into S . Since $(v, \mu_0, \text{Par}_1, \text{Par}_2)$ are all functions of λ , we update them by UPDATE_BY_LAMBDA, in which λ^{inc} denotes the increment to reach the next turning point from the current one. Unlike (v, μ_0) , $(\text{Par}_1, \text{Par}_2)$ has discontinuity at each turning point λ due to the change of support S . They are updated by UPDATE_SHRINK_SUPPORT and UPDATE_EXPAND_SUPPORT depending on whether S is shrunk or expanded. The procedure is repeated until λ crosses 1 and an inner-loop finishes. At the end, Par_2 is recomputed for new $g^{(t+1)}$, which is achieved by DIRECT_UPDATE.

B.3 FIND_LAMBDA

With the help of intermediate variables, we can express (x_S, μ_{S^c}, μ_0) in a compact way.

Theorem 4 *Before any entry of (x_S, μ_{S^c}) hitting 0, it holds that*

$$\mu_0(\lambda) = \mu_0 + \frac{\alpha(\lambda)}{D - \alpha(\lambda) D_g^2} \cdot D_g (D_g \mu_0 - D_{gr}). \quad (23)$$

and

$$v(\lambda) \triangleq \begin{pmatrix} x_S(\lambda) \\ -\mu_{S^c}(\lambda) \end{pmatrix} = v + \frac{\alpha(\lambda)}{D - \alpha(\lambda) D_g^2} \cdot (D_g \mu_0 - D_{gr}) \cdot (D_g \tilde{\eta} - D \eta), \quad (24)$$

Algorithm 3 HOP Algorithm for time-varying A and fixed r

Inputs: Initial matrix $A^{(0)}$, vectors r , matrix-update-vectors $\{g^{(t)}, t = 1, 2, \dots\}$.

Initialization:

$x \leftarrow$ as the optimum corresponding to $A^{(0)}$;
 $S \leftarrow \text{supp}(x)$;
Calculate (x, μ, μ_0) via (10)-(12)
 $v_S \leftarrow x_S, v_{S^c} \leftarrow -\mu_{S^c}$;
Calculate intermediate variables $(\text{Par}_1, \text{Par}_2)$ via (18)-(20) with $g = g^{(1)}$.

Procedure:

```
1: for  $t = 1, 2, \dots$  do
2:    $\lambda \leftarrow 0$ ;
3:   while  $\lambda < 1$  do
4:      $(\lambda^{\text{inc}}, j, S^{\text{new}}) \leftarrow \text{FIND\_LAMBDA}(S, v, \mu_0; \text{Par}_1, \text{Par}_2)$ ;
5:      $\lambda^{\text{inc}} \leftarrow \min\{\lambda^{\text{inc}}, 1 - \lambda\}$ ;
6:      $\lambda \leftarrow \lambda + \lambda^{\text{inc}}$ ;
7:      $(v, \mu_0; \text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_BY\_LAMBDA}(\lambda^{\text{inc}}; v, \mu_0; \text{Par}_1, \text{Par}_2)$ ;
8:     if  $S^{\text{new}} = S \cup \{j\}$  then
9:        $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_EXPAND\_SUPPORT}(\lambda, S, j; r, g^{(t)}, \text{Par}_1, \text{Par}_2)$ ;
10:    else if  $S^{\text{new}} = S \setminus \{j\}$  then
11:       $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_SHRINK\_SUPPORT}(S, j; r, g^{(t)}, \text{Par}_1, \text{Par}_2)$ ;
12:    end if
13:     $S \leftarrow S^{\text{new}}$ .
14:  end while
15:   $\text{Par}_2 \leftarrow \text{DIRECT\_UPDATE}(S, r, g^{(t+1)}; \text{Par}_1, \text{Par}_2)$ ;
16:   $A \leftarrow A + g^{(t)}(g^{(t)})^T$ ;
17:   $x_S^{(t)} \leftarrow x_S, x_{S^c}^{(t)} \leftarrow 0$ .
18: end for
```

Output: $x^{(1)}, x^{(2)}, \dots$.

Proof First we prove (23). By definition,

$$\mu_0(\lambda) = \frac{1 - \mathbf{1}_S^T A_{SS}(\lambda)^{-1} r_S}{\mathbf{1}_S^T A_{SS}(\lambda)^{-1} \mathbf{1}_S} = \frac{1 - \mathbf{1}_S^T M_{SS}(\lambda) r_S}{D(\lambda)}.$$

By Lemma 1,

$$\begin{aligned} -M_{SS}(\lambda) r_S &= -(M_{SS} - \alpha(\lambda) \eta_S \eta_S^T) r_S \\ &= -M_{SS} r_S + \alpha(\lambda) \eta_S \eta_S^T r_S = -M_{SS} r_S - \alpha(\lambda) D_{gr} \eta_S. \end{aligned}$$

Thus the numerator of $\mu_0(\lambda)$ can be written as

$$1 - \mathbf{1}_S^T y_S + \mathbf{1}_S^T (M_{S^c S}^T y_{S^c} - M_{SS} r_S) - \alpha(\lambda) D_{gr} \mathbf{1}_S^T \eta_S = D\mu_0 - \alpha(\lambda) D_{gr} D_g.$$

The denominator of $\mu_0(\lambda)$, by Lemma 1, is formulated as

$$D(\lambda) = D - \alpha(\lambda) D_g^2.$$

Putting the pieces together results in

$$\mu_0(\lambda) = \frac{D\mu_0 - \alpha(\lambda) D_{gr} D_g}{D - \alpha(\lambda) D_g^2} = \mu_0 + \frac{\alpha(\lambda)}{D - \alpha(\lambda) D_g^2} \cdot D_g (D_g \mu_0 - D_{gr}).$$

Plug $\mu_0(\lambda)$ into (10), we obtain that

$$\begin{aligned} x_S(\lambda) &= \mu_0(\lambda) \tilde{\eta}_S(\lambda) + A_{SS}(\lambda)^{-1} r_S \\ &= A_{SS}^{-1} r_S + \alpha(\lambda) D_{gr} \eta_S + (\mu_0(\lambda) - \mu_0) \tilde{\eta}_S(\lambda) + \mu_0 \tilde{\eta}_S(\lambda) \\ &= x_S + \alpha(\lambda) D_{gr} \eta_S + (\mu_0(\lambda) - \mu_0) \tilde{\eta}_S(\lambda) + \mu_0 (\tilde{\eta}_S(\lambda) - \tilde{\eta}_S) \\ &= x_S + \alpha(\lambda) D_{gr} \eta_S + \frac{\alpha(\lambda)}{D - \alpha(\lambda) D_g^2} \cdot D_g (D_g \mu_0 - D_{gr}) \tilde{\eta}_S(\lambda) - \mu_0 \alpha(\lambda) D_g \eta_S \quad [\text{Use Lemma1}] \end{aligned}$$

$$\begin{aligned}
&= x_S + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot D_g(D_g\mu_0 - D_{gr})\tilde{\eta}_S(\lambda) - \alpha(\lambda)(D_g\mu_0 - D_{gr})\eta_S \\
&= x_S + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr}) \cdot (D_g\tilde{\eta}_S(\lambda) - (D - \alpha(\lambda)D_g^2)\eta_S) \\
&= x_S + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr}) \cdot (D_g\tilde{\eta}_S - \alpha(\lambda)D_g^2\eta_S - (D - \alpha(\lambda)D_g^2)\eta_S) \\
&= x_S + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr}) \cdot (D_g\tilde{\eta}_S - D\eta_S).
\end{aligned}$$

Similarly, it follows from (11) that

$$\begin{aligned}
-\mu_{S^c}(\lambda) &= \mu_0(\lambda)\tilde{\eta}_{S^c}(\lambda) + r_{S^c} + M_{S^cS}(\lambda)r_S \\
&= -\mu_{S^c} - \mu_0\tilde{\eta}_{S^c} + (M_{S^cS}(\lambda) - M_{S^cS})r_S + \mu_0(\lambda)\tilde{\eta}_{S^c}(\lambda) \\
&= -\mu_{S^c} - \mu_0\tilde{\eta}_{S^c} + \alpha(\lambda)D_{gr}\eta_{S^c} + \mu_0(\lambda)\tilde{\eta}_{S^c}(\lambda) \\
&= -\mu_{S^c} + \alpha(\lambda)D_{gr}\eta_{S^c} + (\mu_0(\lambda) - \mu_0)\tilde{\eta}_{S^c}(\lambda) + \mu_0(\tilde{\eta}_{S^c}(\lambda) - \tilde{\eta}_{S^c}) \\
&= -\mu_{S^c} + \alpha(\lambda)D_{gr}\eta_{S^c} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot D_g(D_g\mu_0 - D_{gr})\tilde{\eta}_{S^c}(\lambda) - \mu_0\alpha(\lambda)D_g\eta_{S^c} \\
&= -\mu_{S^c} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot D_g(D_g\mu_0 - D_{gr})\tilde{\eta}_{S^c}(\lambda) - \alpha(\lambda)(D_g\mu_0 - D_{gr})\eta_{S^c} \\
&= -\mu_{S^c} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr})(D_g\tilde{\eta}_{S^c}(\lambda) - (D - \alpha(\lambda)D_g^2)\eta_{S^c}) \\
&= -\mu_{S^c} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr})(D_g\tilde{\eta}_{S^c} - \alpha(\lambda)D_g^2\eta_{S^c} - (D - \alpha(\lambda)D_g^2)\eta_{S^c}) \\
&= -\mu_{S^c} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr})(D_g\tilde{\eta}_{S^c} - D\eta_{S^c})
\end{aligned}$$

In sum,

$$\begin{pmatrix} x_S(\lambda) \\ -\mu_{S^c}(\lambda) \end{pmatrix} = \begin{pmatrix} x_S \\ -\mu_{S^c} \end{pmatrix} + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2} \cdot (D_g\mu_0 - D_{gr}) \cdot (D_g\tilde{\eta} - D\eta).$$

■

Theorem 4 indicates that searching for next λ is equivalent to solve n linear equations. In fact, (24) can be abbreviated as

$$v(\lambda) = v + \frac{\alpha(\lambda)}{D - \alpha(\lambda)D_g^2}u = \frac{Dv - (D_g^2v - u)\alpha(\lambda)}{D - \alpha(\lambda)D_g^2},$$

for $u = (D_g\mu_0 - D_{gr}) \cdot (D_g\tilde{\eta} - D\eta)$. Let

$$\alpha = \min_+ \left\{ \frac{Dv_i}{D_g^2v_i - u_i} : i = 1, 2, \dots, n \right\}$$

where $\min_+(\Omega)$ denotes the minimum of all positive numbers contained in set Ω . Then the target λ is the solution of $\alpha(\lambda) = \alpha$, i.e.

$$\lambda = \frac{\alpha}{1 - \alpha D_{gg}}.$$

We should emphasize that the right-handed side might be negative if $\alpha D_{gg} \geq 1$ in which case v never hits 0. Thus, we should set λ to be infinity. The implementation of FIND_LAMBDA is stated in Algorithm 4.

B.4 Variables Update

B.4.1 UPDATE_BY_LAMBDA

Once the next λ has been calculated, all variables can be updated via Lemma 1 and Theorem 4.

Algorithm 4 FIND_LAMBDA

Input: Support S , iterate $v = \begin{pmatrix} x_S \\ -\mu_{S^c} \end{pmatrix}$, μ_0 , intermediate variables $\text{Par}_1, \text{Par}_2$.

Procedure:

```
1:  $u \leftarrow (D_g \mu_0 - D_{gr})(D_g \tilde{\eta} - D\eta)$ ;  
2:  $\alpha \leftarrow \min_+ \left\{ \frac{Dv_i}{D_g^2 v_i - u_i} : i = 1, 2, \dots, n \right\}$ ;  
3:  $j \leftarrow \operatorname{argmin}_+ \left\{ \frac{Dv_i}{D_g^2 v_i - u_i} : i = 1, 2, \dots, n \right\}$ ;  
4: if  $\alpha D_{gg} < 1$  then  
5:    $\lambda^{\text{inc}} \leftarrow \frac{\alpha}{1 - \alpha D_{gg}}$ ;  
6: else  
7:    $\lambda^{\text{inc}} \leftarrow \infty$ ;  
8: end if  
9: if  $j \in S$  then  
10:   $S^{\text{new}} = S \setminus \{j\}$ ;  
11: else  
12:   $S^{\text{new}} = S \cup \{j\}$ .  
13: end if
```

Output: $\lambda^{\text{inc}}, j, S^{\text{new}}$.

Algorithm 5 UPDATE_BY_LAMBDA

Input: Increment λ^{inc} ; iterate $v = \begin{pmatrix} x_S \\ -\mu_{S^c} \end{pmatrix}$, μ_0 ; intermediate variables $\text{Par}_1, \text{Par}_2$.

Procedure:

```
1:  $\alpha_0 \leftarrow \frac{1}{1 + \lambda^{\text{inc}} \cdot D_{gg}}$ ;  
2:  $\alpha \leftarrow \lambda^{\text{inc}} \cdot \alpha_0$ ;  
3:  $\tilde{\alpha} \leftarrow \frac{\alpha}{D - \alpha D_g^2}$ ;  
4:  $v \leftarrow v + \tilde{\alpha} \cdot (D_g \mu_0 - D_{gr})(D_g \tilde{\eta} - D\eta)$ ;  
5:  $\mu_0 \leftarrow \mu_0 + \tilde{\alpha} \cdot D_g(D_g \mu_0 - D_{gr})$ ;  
6:  $D \leftarrow D - \alpha D_g^2$ ;  
7:  $(D_g, D_{gg}, D_{gr}) \leftarrow \alpha_0(D_g, D_{gg}, D_{gr})$ .  
8:  $M_{\cdot, S} \leftarrow M_{\cdot, S} - \alpha \eta \eta_S^T$ ;  
9:  $\tilde{\eta} \leftarrow \tilde{\eta} - \alpha D_g \eta$ ;  
10:  $\eta \leftarrow \alpha_0 \eta$ ;
```

Output: $v, \mu_0, \text{Par}_1, \text{Par}_2$.

B.4.2 UPDATE_EXPAND_SUPPORT

Suppose S is updated to $S \cup \{j\}$ for some $j \in S^c$. Denote \tilde{S} by $S \cup \{j\}$ and we add a superscript $+$ to each variable to denote the value after update. The key tool is the following formula showing the relation between matrix inverses after adding one row and one column.

Proposition 5 Let $\tilde{A}_{jj} = A_{jj} - A_{jS} A_{SS}^{-1} A_{Sj}$,

$$A_{\tilde{S}\tilde{S}}^{-1} = \begin{pmatrix} A_{SS}^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \frac{1}{\tilde{A}_{jj}} \cdot \begin{pmatrix} -A_{SS}^{-1} A_{Sj} \\ 1 \end{pmatrix} (-A_{jS} A_{SS}^{-1} \quad 1).$$

Similar to section 4.1, the key is to update M and other variables are easy to update based on M . Define a class of operators $\{\mathcal{R}_j : j \in \{1, \dots, n\}\}$ such that for matrix $W \in \mathbb{R}^{n \times n}$, $\mathcal{R}_j(W)$ sets the j -th row and j -th column of W to be zero and for vector $z \in \mathbb{R}^{n \times 1}$, $\mathcal{R}_j(z)$ sets the j -th coordinate of z to be zero. One property of \mathcal{R}_j to be used is that For any matrix-vector pair (W, z) ,

$$\mathcal{R}_j(W)z = \mathcal{R}_j(Wz) - z_j \mathcal{R}_j(W_j) \quad (25)$$

where W_j is j -th column of W .

Theorem 6 Let γ and $\tilde{\gamma}$ be two $n \times 1$ vectors with

$$\gamma_{\tilde{S}} = \tilde{\gamma}_{\tilde{S}} = (M_{jS} \quad 1)^T, \quad \gamma_{\tilde{S}^c} = -A_{\tilde{S}^c j} - A_{\tilde{S}^c S} M_{jS}^T, \quad \tilde{\gamma}_{\tilde{S}^c} = 0.$$

Then

$$M^+ = \mathcal{R}_j(M) + \frac{1}{\tilde{A}_{jj}} \cdot \gamma \tilde{\gamma}^T.$$

Proof By definition,

$$M_{\tilde{S}\tilde{S}}^+ = A_{\tilde{S}\tilde{S}}^{-1}, \quad M_{\tilde{S}^c\tilde{S}}^+ = -A_{\tilde{S}^c\tilde{S}} A_{\tilde{S}\tilde{S}}^{-1}.$$

By Proposition 5,

$$M_{\tilde{S}\tilde{S}}^+ = \begin{pmatrix} M_{SS} & 0 \\ 0 & 0 \end{pmatrix} + \frac{1}{\tilde{A}_{jj}} \begin{pmatrix} M_{jS}^T \\ 1 \end{pmatrix} (M_{jS} \quad 1) = \begin{pmatrix} M_{SS} & 0 \\ 0 & 0 \end{pmatrix} + \frac{1}{\tilde{A}_{jj}} \gamma_{\tilde{S}} \tilde{\gamma}_{\tilde{S}}^T,$$

and

$$\begin{aligned} M_{\tilde{S}^c\tilde{S}}^+ &= -(A_{\tilde{S}^c S} \quad A_{\tilde{S}^c j}) \left\{ \begin{pmatrix} M_{SS} & 0 \\ 0 & 0 \end{pmatrix} + \frac{1}{\tilde{A}_{jj}} \begin{pmatrix} M_{jS}^T \\ 1 \end{pmatrix} (M_{jS} \quad 1) \right\} \\ &= (M_{\tilde{S}^c S} \quad 0) + \frac{1}{\tilde{A}_{jj}} \gamma_{\tilde{S}^c} \tilde{\gamma}_{\tilde{S}}^T. \end{aligned}$$

Note that M_{\cdot, \tilde{S}^c} is always a zero matrix by definition, the above results imply that

$$M^+ = \mathcal{R}_j(M) + \frac{1}{\tilde{A}_{jj}} \cdot \gamma \tilde{\gamma}^T.$$

■

The update of other parameters can be derived as a consequence of Theorem 6. Theorem 7 summarizes the results.

Theorem 7 Let $b_{j,S} = -r_{\tilde{S}}^T \gamma_{\tilde{S}}$, then

- $\eta^+ = \mathcal{R}_j(\eta) + \frac{\eta_j}{A_{jj}} \gamma;$
- $\tilde{\eta}^+ = \mathcal{R}_j(\tilde{\eta}) + \frac{\tilde{\eta}_j}{\tilde{A}_{jj}} \gamma;$
- $D^+ = D + \frac{\tilde{\eta}_j^2}{\tilde{A}_{jj}}$
- $D_g^+ = D_g + \frac{\eta_j \tilde{\eta}_j}{\tilde{A}_{jj}};$
- $D_{gg}^+ = D_{gg} + \frac{\eta_j^2}{\tilde{A}_{jj}};$
- $D_{gr}^+ = D_{gr} + \frac{\eta_j b_{j,S}}{\tilde{A}_{jj}};$

Proof Since $M_j = 0$, (25) implies that for any $z \in \mathbb{R}^{n \times 1}$

$$\mathcal{R}_j(Mz) = \mathcal{R}_j(M)z.$$

By definition,

$$\eta = \begin{pmatrix} 0 \\ g_{S^c} \end{pmatrix} + Mg, \quad \tilde{\eta} = \begin{pmatrix} 0 \\ \mathbf{1}_{S^c} \end{pmatrix} + M\mathbf{1}.$$

Also notice that $\tilde{\gamma}_{\tilde{S}}^T g_{\tilde{S}} = g_j + M_{jS}^T g_S = \eta_j$ and $\gamma_{\tilde{S}}^T \mathbf{1}_{\tilde{S}} = 1 + M_{jS} \mathbf{1}_S = \tilde{\eta}_j$, thus,

$$\eta^+ = \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + M^+ g = \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(M)g + \frac{\tilde{\gamma}_{\tilde{S}}^T g}{\tilde{A}_{jj}} \gamma$$

$$\begin{aligned}
&= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(Mg) + \frac{\tilde{\gamma}_{\tilde{S}}^T g_{\tilde{S}}}{\tilde{A}_{jj}} \gamma \\
&= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(\eta) - \mathcal{R}_j \left(\begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} \right) + \frac{\eta_j}{\tilde{A}_{jj}} \gamma \\
&= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(\eta) - \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \frac{\eta_j}{\tilde{A}_{jj}} \gamma \\
&= \mathcal{R}_j(\eta) + \frac{\eta_j}{\tilde{A}_{jj}} \gamma.
\end{aligned}$$

The update of $\tilde{\eta}$ can be obtained by replacing g by $\mathbf{1}$ in the above derivation. The four scalars D, D_g, D_{gg}, D_{gr} can be updated as follows.

$$\begin{aligned}
D^+ &= \mathbf{1}_{\tilde{S}}^T \tilde{\eta}_{\tilde{S}}^+ = \mathbf{1}_{\tilde{S}}^T \left(\mathcal{R}_j(\tilde{\eta})_{\tilde{S}} + \frac{\tilde{\eta}_j}{\tilde{A}_{jj}} \gamma_{\tilde{S}} \right) = D + \frac{\tilde{\eta}_j^2}{\tilde{A}_{jj}}; \\
D_g^+ &= \mathbf{1}_{\tilde{S}}^T \eta_{\tilde{S}}^+ = \mathbf{1}_{\tilde{S}}^T \left(\mathcal{R}_j(\eta)_{\tilde{S}} + \frac{\eta_j}{\tilde{A}_{jj}} \gamma_{\tilde{S}} \right) = D_g + \frac{\tilde{\eta}_j \eta_j}{\tilde{A}_{jj}}; \\
D_{gg}^+ &= g_{\tilde{S}}^T \eta_{\tilde{S}}^+ = g_{\tilde{S}}^T \left(\mathcal{R}_j(\eta)_{\tilde{S}} + \frac{\eta_j}{\tilde{A}_{jj}} \gamma_{\tilde{S}} \right) = D_{gg} + \frac{\eta_j^2}{\tilde{A}_{jj}}; \\
D_{gr}^+ &= -r_{\tilde{S}}^T \eta_{\tilde{S}}^+ = -r_{\tilde{S}}^T \left(\mathcal{R}_j(\eta)_{\tilde{S}} + \frac{\eta_j}{\tilde{A}_{jj}} \gamma_{\tilde{S}} \right) = D_{gr} - \frac{\eta_j}{\tilde{A}_{jj}} r_{\tilde{S}}^T \gamma_{\tilde{S}} \\
&= D_{gr} + \frac{\eta_j b_{j,S}}{\tilde{A}_{jj}}.
\end{aligned}$$

■

The implementation of UPDATE_EXPAND_SUPPORT is summarized in Algorithm 6. Note that both \tilde{A}_{jj} and $\gamma_{\tilde{S}^c}$ depend on λ and it is easy to see that

$$\begin{aligned}
\tilde{A}_{jj}(\lambda) &= A_{jj} + \lambda g_j^2 + M_{jS}(A_{Sj} + \lambda g_j g_S) = A_{jj} + M_{jS} A_{Sj} + \lambda g_j \eta_j \\
\gamma_{\tilde{S}^c}(\lambda) &\leftarrow -(A_{\tilde{S}^c j} + \lambda g_{\tilde{S}^c} g_j) - (A_{\tilde{S}^c S} + \lambda g_{\tilde{S}^c} g_S^T) M_{jS}^T = -A_{\tilde{S}^c j} - A_{\tilde{S}^c S} M_{jS}^T - \lambda \eta_j g_{\tilde{S}^c}.
\end{aligned}$$

B.4.3 UPDATE_SHRINK_SUPPORT

Suppose S is updated to $S \setminus \{j\}$ for some $j \in S^c$. Denote \tilde{S} by $S \setminus \{j\}$ and we add a superscript $-$ to each variable to denote the value after update. Similar to last subsection, we start from deriving M^- and apply the result to calculate other variables.

Theorem 8 Let β and $\tilde{\beta}$ be two $n \times 1$ vectors with

$$\beta_{\tilde{S}} = \tilde{\beta}_{\tilde{S}} = M_{\tilde{S}j}, \quad \beta_{\tilde{S}^c} = \begin{pmatrix} -1 \\ M_{S^c j} \end{pmatrix} \quad \tilde{\beta}_{\tilde{S}^c} = 0.$$

Then

$$M^- = \mathcal{R}_j(M) - \frac{1}{M_{jj}} \cdot \beta \tilde{\beta}^T.$$

Proof By definition,

$$\begin{pmatrix} M_{\tilde{S}\tilde{S}} & M_{\tilde{S}j} \\ M_{j\tilde{S}} & M_{jj} \end{pmatrix} = A_{SS}^{-1} = \begin{pmatrix} A_{\tilde{S}\tilde{S}} & A_{\tilde{S}j} \\ A_{j\tilde{S}} & A_{jj} \end{pmatrix}^{-1}.$$

Then Proposition 5 implies that

$$\begin{pmatrix} M_{\tilde{S}\tilde{S}} & M_{\tilde{S}j} \\ M_{j\tilde{S}} & M_{jj} \end{pmatrix} = \begin{pmatrix} A_{\tilde{S}\tilde{S}}^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \frac{1}{A_{jj} - A_{j\tilde{S}} A_{\tilde{S}\tilde{S}}^{-1} A_{\tilde{S}j}} \cdot \begin{pmatrix} -A_{\tilde{S}\tilde{S}}^{-1} A_{\tilde{S}j} \\ 1 \end{pmatrix} \begin{pmatrix} -A_{j\tilde{S}} A_{\tilde{S}\tilde{S}}^{-1} & 1 \end{pmatrix}. \quad (26)$$

Algorithm 6 UPDATE_EXPAND_SUPPORT

Inputs: Current λ , original support S , new index j , matrix A , vectors y, r, g , intermediate variables $\text{Par}_1, \text{Par}_2$.

Procedure:

- 1: $\tilde{A}_{jj} \leftarrow A_{jj} + M_{jS}A_{Sj} + \lambda g_j \eta_j$;
- 2: $\gamma_{\tilde{S}} \leftarrow (M_{jS}, 1)^T, \gamma_{\tilde{S}^c} \leftarrow -A_{\tilde{S}^c j} - A_{\tilde{S}^c S} M_{jS}^T - \lambda \eta_j g_{\tilde{S}^c}$;
- 3: $\tilde{\gamma}_{\tilde{S}} \leftarrow (M_{jS}, 1)^T, \tilde{\gamma}_{\tilde{S}^c} \leftarrow 0$;
- 4: $b \leftarrow -r_{\tilde{S}}^T \gamma_{\tilde{S}}$;
- 5: $D \leftarrow D + \frac{\eta_j^2}{A_{jj}}$;
- 6: $D_g \leftarrow D_g + \frac{\eta_j \tilde{\eta}_i}{A_{jj}}$;
- 7: $D_{gg} \leftarrow D_{gg} + \frac{\eta_j^2}{A_{jj}}$;
- 8: $D_{gr} \leftarrow D_{gr} + \frac{\eta_j b}{A_{jj}}$;
- 9: $M_{\cdot, \tilde{S}} \leftarrow \mathcal{R}_j(M_{\cdot, \tilde{S}}) + \frac{1}{A_{jj}} \gamma \tilde{\gamma}_{\tilde{S}}^T$;
- 10: $\eta \leftarrow \mathcal{R}_j(\eta) + \frac{\eta_j}{A_{jj}} \gamma$.
- 11: $\tilde{\eta} \leftarrow \mathcal{R}_j(\tilde{\eta}) + \frac{\tilde{\eta}_j}{A_{jj}} \gamma$;

Output: $\text{Par}_1, \text{Par}_2$.

This entails that

$$A_{\tilde{S}\tilde{S}}^{-1} = M_{\tilde{S}\tilde{S}} - \frac{M_{\tilde{S}j}M_{j\tilde{S}}}{M_{jj}}, \quad -A_{j\tilde{S}}A_{\tilde{S}\tilde{S}}^{-1} = \frac{M_{j\tilde{S}}}{M_{jj}}. \quad (27)$$

On the other hand,

$$\begin{aligned} (M_{S^c\tilde{S}} \quad M_{S^c j}) &= -A_{S^c S}A_{\tilde{S}\tilde{S}}^{-1} = -(A_{S^c\tilde{S}} \quad A_{S^c j}) \begin{pmatrix} M_{\tilde{S}\tilde{S}} & M_{\tilde{S}j} \\ M_{j\tilde{S}} & M_{jj} \end{pmatrix} \\ &= -(A_{S^c\tilde{S}}M_{\tilde{S}\tilde{S}} + A_{S^c j}M_{j\tilde{S}} \quad A_{S^c\tilde{S}}M_{\tilde{S}j} + A_{S^c j}M_{jj}). \end{aligned} \quad (28)$$

It follows from (26), (27) and (28) that

$$\begin{aligned} -A_{S^c\tilde{S}}A_{\tilde{S}\tilde{S}}^{-1} &= -A_{S^c\tilde{S}} \left(M_{\tilde{S}\tilde{S}} - \frac{M_{\tilde{S}j}M_{j\tilde{S}}}{M_{jj}} \right) \\ &= M_{S^c\tilde{S}} + A_{S^c j}M_{j\tilde{S}} + \frac{A_{S^c\tilde{S}}M_{\tilde{S}j}M_{j\tilde{S}}}{M_{jj}} \\ &= M_{S^c\tilde{S}} + (A_{S^c j}M_{jj} + A_{S^c\tilde{S}}M_{\tilde{S}j}) \frac{M_{j\tilde{S}}}{M_{jj}} \\ &= M_{S^c\tilde{S}} - \frac{M_{S^c j}M_{j\tilde{S}}}{M_{jj}}. \end{aligned} \quad (29)$$

Putting (27) and (28) together, we obtain that

$$M_{\cdot, \tilde{S}}^{-} = \begin{pmatrix} A_{\tilde{S}\tilde{S}}^{-1} \\ -A_{\tilde{S}^c\tilde{S}}A_{\tilde{S}\tilde{S}}^{-1} \end{pmatrix} = \begin{pmatrix} A_{\tilde{S}\tilde{S}}^{-1} \\ -A_{j\tilde{S}}A_{\tilde{S}\tilde{S}}^{-1} \\ -A_{S^c\tilde{S}}A_{\tilde{S}\tilde{S}}^{-1} \end{pmatrix} = \mathcal{R}_j(M)_{\cdot, \tilde{S}} - \frac{1}{M_{jj}} \cdot \beta \tilde{\beta}_{\tilde{S}}^T.$$

Since M_{\cdot, \tilde{S}^c} is a zero matrix,

$$M^{-} = \mathcal{R}_j(M) - \frac{1}{M_{jj}} \cdot \beta \tilde{\beta}^T.$$

■

Theorem 9 Let $\tilde{b}_{j,S} = -r_{\tilde{S}}^T \beta_{\tilde{S}} - r_j M_{jj}$, then

- $\eta^{-} = \mathcal{R}_j(\eta) - \frac{\eta_j}{M_{jj}} \beta$;

Algorithm 7 UPDATE_SHRINK_SUPPORT

Inputs: Original support S , new index j , matrix A , vector y, r, g , intermediate variables $\text{Par}_1, \text{Par}_2$.

Procedure:

- 1: $\beta_{\tilde{S}} \leftarrow M_{j\tilde{S}}^T, \beta_{\tilde{S}^c} \leftarrow \begin{pmatrix} -1 \\ M_{\tilde{S}^c j} \end{pmatrix}, \tilde{\beta}_{\tilde{S}} \leftarrow M_{j\tilde{S}}^T, \tilde{\beta}_{\tilde{S}^c} \leftarrow 0;$
- 2: $\tilde{b} \leftarrow -r_{\tilde{S}}^T \beta_{\tilde{S}} - r_j M_{jj};$
- 3: $D \leftarrow D - \frac{\tilde{\eta}_j^2}{M_{jj}};$
- 4: $D_g \leftarrow D_g - \frac{\eta_j \tilde{\eta}_j}{M_{jj}};$
- 5: $D_{gg} \leftarrow D_{gg} - \frac{\eta_j^2}{M_{jj}};$
- 6: $D_{gr} \leftarrow D_{gr} - \frac{\eta_j \tilde{b}}{M_{jj}};$
- 7: $M_{\cdot, \tilde{S}} \leftarrow \mathcal{R}_j(M_{\cdot, \tilde{S}}) - \frac{1}{M_{jj}} \beta \tilde{\beta}_{\tilde{S}}^T, M_{\cdot, j} \leftarrow 0;$
- 8: $\eta \leftarrow \mathcal{R}_j(\eta) - \frac{\eta_j}{M_{jj}} \beta;$
- 9: $\tilde{\eta} \leftarrow \mathcal{R}_j(\tilde{\eta}) - \frac{\tilde{\eta}_j}{M_{jj}} \beta.$

Output: $\text{Par}_1, \text{Par}_2$.

- $\tilde{\eta}^- = \mathcal{R}_j(\tilde{\eta}) - \frac{\tilde{\eta}_j}{M_{jj}} \beta;$
- $D^- = D - \frac{\tilde{\eta}_j^2}{M_{jj}};$
- $D_g^- = D_g - \frac{\eta_j \tilde{\eta}_j}{M_{jj}};$
- $D_{gg}^- = D_{gg} - \frac{\eta_j^2}{M_{jj}};$
- $D_{gr}^- = D_{gr} - \frac{\eta_j \tilde{b}_{j, S}}{M_{jj}}.$

Proof By (25),

$$\mathcal{R}_j(M)g = \mathcal{R}_j(Mg) - g_j \mathcal{R}_j(M_{\cdot, j})$$

Let e_j is the j -th basis vector with j -th entry equal to 1 and all other entries equal to 0. Then

$$\begin{aligned} \eta^- &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + M^- g = \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(M)g - \frac{\tilde{\beta}_{\tilde{S}}^T g}{M_{jj}} \beta \\ &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(Mg) - g_j \mathcal{R}_j(M_{\cdot, j}) - \frac{\tilde{\beta}_{\tilde{S}}^T g_{\tilde{S}}}{M_{jj}} \beta \\ &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(Mg) - g_j e_j - g_j \beta - \frac{\tilde{\beta}_{\tilde{S}}^T g_{\tilde{S}}}{M_{jj}} \beta \\ &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(Mg) - g_j e_j - \frac{M_{jS} g_S}{M_{jj}} \beta \\ &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(\eta) - \mathcal{R}_j \left(\begin{pmatrix} 0 \\ g_{S^c} \end{pmatrix} \right) - g_j e_j - \frac{\eta_j}{M_{jj}} \beta \\ &= \begin{pmatrix} 0 \\ g_{\tilde{S}^c} \end{pmatrix} + \mathcal{R}_j(\eta) - \begin{pmatrix} 0 \\ g_{S^c} \end{pmatrix} - g_j e_j - \frac{\eta_j}{M_{jj}} \beta \\ &= \mathcal{R}_j(\eta) - \frac{\eta_j}{M_{jj}} \beta. \end{aligned}$$

Substitute g by $\mathbf{1}$, we obtain the update for $\tilde{\eta}$. Together with (27) and the fact that $j \in S$, we obtain that

$$D^- = \mathbf{1}_S^T \tilde{\eta}_{\tilde{S}}^- = \mathbf{1}_S^T \left(\tilde{\eta}_{\tilde{S}} - \frac{\tilde{\eta}_j}{M_{jj}} \beta_{\tilde{S}} \right) = D - \tilde{\eta}_j - \frac{\tilde{\eta}_j}{M_{jj}} (\mathbf{1}_S^T M_{\tilde{S}j})$$

$$\begin{aligned}
&= D - \frac{\tilde{\eta}_j}{M_{jj}}(M_{jj} + \mathbf{1}_{\tilde{S}}^T M_{\tilde{S}j}) = D - \frac{\tilde{\eta}_j^2}{M_{jj}}; \\
D_g^- &= \mathbf{1}_{\tilde{S}}^T \eta_{\tilde{S}}^- = \mathbf{1}_{\tilde{S}}^T \left(\eta_{\tilde{S}} - \frac{\eta_j}{M_{jj}} \beta_{\tilde{S}} \right) = D_g - \eta_j - \frac{\eta_j}{M_{jj}} (\mathbf{1}_{\tilde{S}}^T M_{\tilde{S}j}) \\
&= D_g - \frac{\eta_j}{M_{jj}} (M_{jj} + \mathbf{1}_{\tilde{S}}^T M_{\tilde{S}j}) = D_g - \frac{\eta_j \tilde{\eta}_j}{M_{jj}}; \\
D_{gg}^- &= g_{\tilde{S}}^T \eta_{\tilde{S}}^- = g_{\tilde{S}}^T \left(\eta_{\tilde{S}} - \frac{\eta_j}{M_{jj}} \beta_{\tilde{S}} \right) = D_{gg} - g_j \eta_j - \frac{\eta_j}{M_{jj}} (g_{\tilde{S}}^T M_{\tilde{S}j}) \\
&= D_{gg} - \frac{\eta_j}{M_{jj}} (g_j M_{jj} + g_{\tilde{S}}^T M_{\tilde{S}j}) = D_{gg} - \frac{\eta_j^2}{M_{jj}}; \\
D_{gr}^- &= -r_{\tilde{S}}^T \eta_{\tilde{S}}^- = -r_{\tilde{S}}^T \left(\eta_{\tilde{S}} - \frac{\eta_j}{M_{jj}} \beta_{\tilde{S}} \right) = -r_{\tilde{S}}^T \eta_S + r_j \eta_j + \frac{\eta_j}{M_{jj}} r_{\tilde{S}}^T \beta_{\tilde{S}} \\
&= D_{gr} - \frac{\eta_j \tilde{b}_{j,S}}{M_{jj}}.
\end{aligned}$$

■

The implementation of UPDATE_SHRINK_SUPPORT is summarized in Algorithm 7.

B.4.4 DIRECT_UPDATE

At the beginning of each time t , we need to recompute $\text{Par}_2 = \{\eta, D_g, D_{gg}, D_{gr}\}$. The implementation is summarized in Algorithm 8.

Algorithm 8 DIRECT_UPDATE

Inputs: Support S , vector y, r, g , intermediate variables $\text{Par}_1, \text{Par}_2$.

Procedure:

- 1: $\eta_S \leftarrow M_{SS} g_S, \eta_{S^c} \leftarrow g_{S^c} + M_{S^c S} g_S;$
- 2: $D_g \leftarrow \mathbf{1}_S^T \eta_S;$
- 3: $D_{gg} \leftarrow \eta_S^T g_S;$
- 4: $D_{gr} \leftarrow -\eta_S^T r_S.$

Output: Par_2 .

B.5 Update of A

As will be shown in next subsection, the complexities of all above sub-routines are at most $O(ns)$ where $s = |S|$. However, the complexity of line 16 in Algorithm 3 is $O(n^2)$ which might dominate when the solution is sparse and the number of turning points is small. Fortunately, UPDATE_EXPAND_SUPPORT is the only sub-routine which extracts information from A . In fact, in line 1 and line 2,

$$\begin{pmatrix} \tilde{A}_{jj} \\ \gamma_{\tilde{S}^c} \end{pmatrix} = \begin{pmatrix} A_{jj} + M_{jS} A_{Sj} \\ -A_{\tilde{S}^c j} - A_{\tilde{S}^c S} M_{jS}^T \end{pmatrix} + \lambda \eta_j \begin{pmatrix} g_j \\ g_{\tilde{S}^c} \end{pmatrix}.$$

This only requires the j -th column of A . Let S_* be the union of all supports appeared in Algorithm 3. Suppose we know S_* apriori, we can only update the columns of A with indices in S_* . In other words, we update A_{\cdot, S_*} by $A_{\cdot, S_*} + \lambda g g_{S_*}^T$ at the beginning of each step and hence the complexity is reduced to $O(n|S_*|)$.

Although agnostic to S_* in reality, we can initialize it by $\text{supp}(x_k)$ for some positive k , e.g. $k = 1$, and keep track it by adding index into S_* once the index is not included in S_* . Once a new index j is detected, we update j -th column of A by using all previous $g^{(t)}$. The implementation is stated in Algorithm 9.

Algorithm 9 HOP Algorithm for time-varying A and fixed r with late update of A

Inputs: Initial matrix $A^{(0)}$, vectors r , matrix-update-vectors $\{g^{(t)}, t = 1, 2, \dots\}$.

Initialization:

$x \leftarrow$ as the optimum corresponding to $A^{(0)}$;
 $S \leftarrow \text{supp}(x)$, $S_* \leftarrow S$;
Calculate (x, μ, μ_0) via (10)-(12)
 $v_S \leftarrow x_S$, $v_{S^c} \leftarrow -\mu_{S^c}$;
Calculate intermediate variables $(\text{Par}_1, \text{Par}_2)$ via (18)-(20) based on $g^{(1)}$.

Procedure:

```

1: for  $t = 1, 2, \dots$  do
2:    $\lambda \leftarrow 0$ ;
3:   while  $\lambda < 1$  do
4:      $(\lambda^{\text{inc}}, j, S^{\text{new}}) \leftarrow \text{FIND\_LAMBDA}(S, v; \text{Par}_1, \text{Par}_2)$ ;
5:      $\lambda^{\text{inc}} \leftarrow \min\{\lambda^{\text{inc}}, 1 - \lambda\}$ ;
6:      $\lambda \leftarrow \lambda + \lambda^{\text{inc}}$ ;
7:      $(v, \mu_0; \text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_BY\_LAMBDA}(\lambda^{\text{inc}}; v, \mu_0; \text{Par}_1, \text{Par}_2)$ ;
8:     if  $S^{\text{new}} = S \cup \{j\}$  then
9:        $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_EXPAND\_SUPPORT}(\lambda, S, j; A, r, g^{(t)}; \text{Par}_1, \text{Par}_2)$ ;
10:      if  $j \notin S_*$  then
11:         $G \leftarrow (g^{(1)}, \dots, g^{(t-1)})$ ;
12:         $A_{:,j} \leftarrow A_{:,j} + GG_{j,:}^T$ ;
13:         $S_* = S_* \cup \{j\}$ ;
14:      end if
15:    else if  $S^{\text{new}} = S \setminus \{j\}$  then
16:       $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_SHRINK\_SUPPORT}(S, j; r, g^{(t)}; \text{Par}_1, \text{Par}_2)$ ;
17:    end if
18:     $S \leftarrow S^{\text{new}}$ .
19:  end while
20:   $\text{Par}_2 \leftarrow \text{DIRECT\_UPDATE}(S, r, g^{(t+1)}; \text{Par}_1, \text{Par}_2)$ ;
21:   $A_{:,S_*} \leftarrow A_{:,S_*} + g^{(t)}(g_{S_*}^{(t)})^T$ ;
22:   $x_S^{(t)} \leftarrow x_S$ ,  $x_{S^c}^{(t)} \leftarrow 0$ .
23: end for

```

Output: $x^{(1)}, x^{(2)}, \dots$.

B.6 Complexity Analysis

In this subsection, we analyze the complexity of the algorithm. We distinguish four types of computation, namely the matrix-vector product, the outer-product of two vectors, the inner-product of two vectors and the vector-scalar product. Denote by $W \in \mathbb{R}^{m \times p}$, $(z, \tilde{z}) \in \mathbb{R}^p \times \mathbb{R}^q$ and $a \in \mathbb{R}$ the generic matrix, vector and scalar respectively. As a convention, the complexity is defined as the number of scalar-scalar multiplications. The addition is omitted here for simplicity. Note that the complexities of Wz , $z\tilde{z}^T$, $z^T z$ and az are mp , pq , p and p , respectively. The results for a single step are summarized in Table 5 where $s_* = |S_*|$ is the size of S_* at the final round. We should emphasize that our complexity analysis is more accurate than big-O analysis in that it reveals the exact constant.

Table 5: Computation complexity of each sub-routine in Algorithm 3.

	(Wz) -type	$(z\tilde{z}^T)$ -type	$(z^T z)$ -type	(az) -type
FIND_LAMBDA	0	0	n	$2n$
UPDATE_BY_LAMBDA	0	ns	0	$4n$
UPDATE_EXPAND_SUPPORT	$s(n - s - 1)$	$n(s + 1)$	n	$2n$
UPDATE_SHRINK_SUPPORT	0	$n(s - 1)$	n	$2n$
DIRECT_UPDATE	ns	0	$n + s$	0
Update ⁶ of A .	0	ns_*	0	0

For given t , denote k_A^+ by the number of turning points which add elements to S and k_A^- by the number of turning points which delete elements from S . Let $k_A = k_A^+ + k_A^-$ be the total number of turning points and s be the maximum size of S in the iteration. Then the complexity of HOP for a single time t is at most

$$ns(3k_A^+ + 2k_A^-) + n(12k_A^+ + 10k_A^-) + O(k_A),$$

Therefore, the complexity at time t is at most

$$C_{1t} \leq ns_* + ns(3k_A^+ + 2k_A^- + 1) + n(12k_A^+ + 10k_A^- + 2) + O(k_A) \leq ns_* + ns(3k_A + 1) + n(12k_A + 2) + O(k_A).$$

C Implementation of HOP Algorithm With Time-Varying r and Fixed A

C.1 Intermediate Variables

Similar to Appendix B, we define $\text{Par}_1 = \{M, \tilde{\eta}, D\}$ where the parameters are defined in (18)-(20). Moreover, we define a vector ξ such that

$$\xi_S = -A_{SS}^{-1}\ell_S, \quad \xi_{S^c} = -\ell_{S^c} + A_{S^cS}A_{SS}^{-1}\ell_S,$$

and a scalar D_ℓ as

$$D_\ell = \mathbf{1}_S^T \xi_S.$$

We write Par_3 for $\{\xi, D_\ell\}$ for convenience.

C.2 Implementation

Algorithm 10 describes the full implementation in this case and the sub-routines will be discussed separately in following subsections.

C.3 FIND_UTILDE_LAMBDA

Define v as in (22). Then Theorem 1 in the main text implies that

$$v(\underline{\lambda}) = v(0) - \left(\xi - \frac{D_\ell}{D} \tilde{\eta} \right) \underline{\lambda}, \quad \mu_0(\underline{\lambda}) = \mu_0 + \frac{D_\ell}{D} \underline{\lambda}.$$

Thus, searching for $\underline{\lambda}$ is equivalent to solving simple linear equations; see Algorithm 11.

C.4 Variables Update

C.4.1 UPDATE_BY_UTILDE_LAMBDA

Note that all intermediate variables are not affected by $\underline{\lambda}$, we only need to update v and μ_0 accordingly.

C.4.2 UPDATE_UTILDE_EXPAND_SUPPORT

Since M is exactly the same as that defined in Appendix B, we can directly apply Theorem 6 to obtain an update of M and the updates of other parameters as a consequence.

Theorem 10 *Let γ and $\tilde{\gamma}$ be defined in Theorem 6, i.e.*

$$\gamma_{\tilde{S}} = \tilde{\gamma}_{\tilde{S}} = (M_{jS} \quad 1)^T, \quad \gamma_{\tilde{S}^c} = -A_{\tilde{S}^c j} - A_{\tilde{S}^c S} M_{jS}^T, \quad \tilde{\gamma}_{\tilde{S}^c} = 0,$$

then

- $M^+ = \mathcal{R}_j(M) + \frac{1}{A_{jj}} \gamma \tilde{\gamma}^T;$
- $\tilde{\eta}^+ = \mathcal{R}_j(\tilde{\eta}) + \frac{\tilde{\eta}_j}{A_{jj}} \gamma;$

⁶ S_* might be updated, in which step nt -computations are involved in line 17 of Algorithm 9. However, on average, ns_* computations are involved since s_* represents the size of S_* at the last round.

Algorithm 10 HOP Algorithm for constant A, y and time-varying r

Inputs: Initial matrix A , vector $r^{(0)}$, vector-update-vector $\{\ell^{(t)} = r^{(t)} - r^{(t-1)} : t = 1, 2, \dots\}$.

Initialization:

$x \leftarrow$ as the optimum corresponding to $r^{(0)}$.
 $S \leftarrow \text{supp}(x)$;
Calculate (x, μ, μ_0) via (10)-(12)
 $v_S \leftarrow x_S, v_{S^c} \leftarrow -\mu_{S^c}$;
Calculate intermediate variables $\text{Par}_1, \text{Par}_3$ via (18)-(20) based on $\ell^{(1)}$.

Procedure:

```
1: for  $t = 1, 2, \dots$  do
2:    $\underline{\lambda} \leftarrow 0$ ;
3:   while  $\underline{\lambda} < 1$  do
4:      $(\underline{\lambda}^{\text{inc}}, j, S^{\text{new}}) \leftarrow \text{FIND\_UTILDE\_LAMBDA}(S, v; \text{Par}_1, \text{Par}_3)$ ;
5:      $\underline{\lambda}^{\text{inc}} \leftarrow \min\{\underline{\lambda}^{\text{inc}}, 1 - \underline{\lambda}\}$ ;
6:      $(v, \mu_0; \text{Par}_1, \text{Par}_3) \leftarrow \text{UPDATE\_BY\_UTILDE\_LAMBDA}(\underline{\lambda}^{\text{inc}}; v, \mu_0, \text{Par}_1, \text{Par}_3)$ ;
7:     if  $S^{\text{new}} = S \cup \{j\}$  then
8:        $(\text{Par}_1, \text{Par}_3) \leftarrow \text{UPDATE\_UTILDE\_EXPAND\_SUPPORT}(S, j, A, \ell; \text{Par}_1, \text{Par}_3)$ ;
9:     else if  $S^{\text{new}} = S \setminus \{j\}$  then
10:       $(\text{Par}_1, \text{Par}_3) \leftarrow \text{UPDATE\_UTILDE\_SHRINK\_SUPPORT}(S, j, \ell; \text{Par}_1, \text{Par}_3)$ ;
11:    end if
12:     $S \leftarrow S^{\text{new}}$ ;
13:     $\underline{\lambda} \leftarrow \underline{\lambda} + \underline{\lambda}^{\text{inc}}$ .
14:  end while
15:   $\text{Par}_3 \leftarrow \text{DIRECT\_UTILDE\_UPDATE}(S, \text{Par}_1, \ell^{(t+1)})$ ;
16:   $x_S^{(t)} \leftarrow x_S, x_{S^c}^{(t)} \leftarrow 0$ .
17: end for
```

Output: $x^{(1)}, x^{(2)}, \dots$.

Algorithm 11 FIND_UTILDE_LAMBDA

Input: Support S , iterate $v = \begin{pmatrix} x_S \\ -\mu_{S^c} \end{pmatrix}$, intermediate variables $\text{Par}_1, \text{Par}_3$.

Procedure:

```
1:  $\underline{\lambda}^{\text{inc}} \leftarrow \min_+ \left\{ \frac{v_i}{\xi_i - \frac{D_\ell}{D} \tilde{\eta}_i} : i = 1, 2, \dots, n \right\}$ ;
2:  $j \leftarrow \text{argmin}_+ \left\{ \frac{v_i}{\xi_i - \frac{D_\ell}{D} \tilde{\eta}_i} : i = 1, 2, \dots, n \right\}$ ;
3: if  $j \in S$  then
4:    $S^{\text{new}} = S \setminus \{j\}$ ;
5: else
6:    $S^{\text{new}} = S \cup \{j\}$ .
7: end if
```

Output: $\underline{\lambda}^{\text{inc}}, j, S^{\text{new}}$.

Algorithm 12 UPDATE_BY_UTILDE_LAMBDA

Input: Increment $\underline{\lambda}^{\text{inc}}$; iterate $v = \begin{pmatrix} x_S \\ -\mu_{S^c} \end{pmatrix}$, μ_0 ; intermediate variables $\text{Par}_1, \text{Par}_3$.

Procedure:

```
1:  $v \leftarrow v - \left( \xi - \frac{D_\ell}{D} \tilde{\eta} \right) \underline{\lambda}^{\text{inc}}$ ;
2:  $\mu_0 \leftarrow \mu_0 + \frac{D_\ell}{D} \underline{\lambda}^{\text{inc}}$ .
```

Output: v, μ_0 .

- $D^+ = D + \frac{\tilde{\eta}_j^2}{A_{jj}};$
- $\xi^+ = \mathcal{R}_j(\xi) + \frac{\xi_j}{A_{jj}}\gamma;$
- $D_\ell^+ = D_\ell + \frac{\xi_j \tilde{\eta}_j}{A_{jj}}.$

Proof The update of M , $\tilde{\eta}$ and D has been proved in Theorem 6. For any subset S , let I_S denote the matrix with j -th diagonal element equal to 1 for any $j \in S$ and all other elements equal to 0. Then ξ and ξ^+ can be rewritten as

$$\xi = -(M + I_{S^c})\ell, \quad \xi^+ = -(M^+ + I_{\tilde{S}^c})\ell.$$

Note that $I_{S^c} - I_{\tilde{S}^c} = e_j e_j^T$ where e_j is the j -th basis vector, then we have

$$\begin{aligned} \xi^+ - \xi &= (M - M^+ + e_j e_j^T)\ell = \left(M - \mathcal{R}_j(M) - \frac{1}{\tilde{A}_{jj}}\gamma\tilde{\gamma}^T + e_j e_j^T \right)\ell = -\frac{\tilde{\gamma}^T \ell}{\tilde{A}_{jj}}\gamma + (\ell_j - M_{jS}\ell_S)e_j \\ \implies \xi^+ &= \xi - \frac{\tilde{\gamma}^T \ell}{\tilde{A}_{jj}}\gamma + (\ell_j - M_{jS}\ell_S)e_j = \mathcal{R}_j(\xi) - \frac{\tilde{\gamma}^T \ell}{\tilde{A}_{jj}}\gamma. \end{aligned}$$

Note that $\tilde{\gamma}^T \ell = \ell_j + M_{jS}\ell_S = -\xi_j$, we obtain that

$$\xi^+ = \mathcal{R}_j(\xi) + \frac{\xi_j}{\tilde{A}_{jj}}\gamma.$$

For D_ℓ^+ , we have

$$D_\ell^+ = \mathbf{1}_{\tilde{S}}^T \xi^+ = D_\ell + \frac{\xi_j}{\tilde{A}_{jj}} \cdot \mathbf{1}_{\tilde{S}}^T \gamma_{\tilde{S}} = D_\ell + \frac{\xi_j \tilde{\eta}_j}{\tilde{A}_{jj}}.$$

■

The implementation of UPDATE_TILDE_EXPAND_SUPPORT is summarized in Algorithm 13.

Algorithm 13 UPDATE_UTILDE_EXPAND_SUPPORT

Inputs: Original support S , new index j , matrix A , vector ℓ , intermediate variables $\text{Par}_1, \text{Par}_3$.

Procedure:

- 1: $\tilde{A}_{jj} \leftarrow A_{jj} + M_{jS}A_{Sj};$
- 2: $\gamma_{\tilde{S}} \leftarrow (M_{jS}, 1)^T, \gamma_{\tilde{S}^c} \leftarrow -A_{\tilde{S}^c j} - A_{\tilde{S}^c S}M_{jS}^T;$
- 3: $\tilde{\gamma}_{\tilde{S}} \leftarrow (M_{jS}, 1)^T, \tilde{\gamma}_{\tilde{S}^c} \leftarrow 0;$
- 4: $D \leftarrow D + \frac{\tilde{\eta}_j^2}{\tilde{A}_{jj}};$
- 5: $D_\ell \leftarrow D_\ell + \frac{\xi_j \tilde{\eta}_j}{\tilde{A}_{jj}};$
- 6: $\xi \leftarrow \mathcal{R}_j(\xi) + \frac{\xi_j}{\tilde{A}_{jj}}\gamma;$
- 7: $\tilde{\eta} \leftarrow \mathcal{R}_j(\tilde{\eta}) + \frac{\tilde{\eta}_j}{\tilde{A}_{jj}}\gamma;$
- 8: $M_{\cdot, \tilde{S}} \leftarrow \mathcal{R}_j(M_{\cdot, \tilde{S}}) + \frac{1}{\tilde{A}_{jj}}\gamma\tilde{\gamma}_{\tilde{S}}^T.$

Output: $\text{Par}_1, \text{Par}_3$.

C.4.3 UPDATE_UTILDE_SHRINK_SUPPORT

Since M is exactly the same as in Appendix B, we can directly apply Theorem 8 to obtain an update of M and the updates of other parameters as a consequence.

Theorem 11 Let β and $\tilde{\beta}$ be defined in Theorem 8, i.e.

$$\beta_{\tilde{S}} = \tilde{\beta}_{\tilde{S}} = M_{\tilde{S}j}, \quad \beta_{\tilde{S}^c} = \begin{pmatrix} -1 \\ M_{\tilde{S}^c j} \end{pmatrix} \quad \tilde{\beta}_{\tilde{S}^c} = 0,$$

then

- $M^- = \mathcal{R}_j(M) - \frac{1}{M_{jj}} \cdot \beta \tilde{\beta}^T$;
- $\tilde{\eta}^- = \mathcal{R}_j(\tilde{\eta}) - \frac{\tilde{\eta}_j}{M_{jj}} \beta$;
- $D^- = D - \frac{\tilde{\eta}_j^2}{M_{jj}}$;
- $\xi^- = \mathcal{R}_j(\xi) - \frac{\xi_j}{M_{jj}} \beta$;
- $D_\ell^- = D_\ell - \frac{\xi_j \tilde{\eta}_j}{M_{jj}}$.

Proof The update of M , $\tilde{\eta}$ and D has been proved in Theorem 8 and Theorem 9. For any subset S , let I_S denote the matrix with j -th diagonal element equal to 1 for any $j \in S$ and all other elements equal to 0. Then ξ and ξ^- can be rewritten as

$$\xi = -(M + I_{S^c}) \ell, \quad \xi^- = -(M^- + I_{\tilde{S}^c}) \ell.$$

Note that $I_{\tilde{S}^c} - I_{S^c} = e_j e_j^T$ where e_j is the j -th basis vector, then we have

$$\begin{aligned} \xi^- - \xi &= (M - M^- - e_j e_j^T) \ell = \left(\frac{1}{M_{jj}} \beta \tilde{\beta}^T + M - \mathcal{R}_j(M) - e_j e_j^T \right) \ell \\ &= \frac{\beta_{\tilde{S}}^T \ell_{\tilde{S}}}{M_{jj}} \beta + (M - \mathcal{R}_j(M) - e_j e_j^T) \ell \triangleq \frac{\beta_{\tilde{S}}^T \ell_{\tilde{S}}}{M_{jj}} \delta + \tilde{\xi}. \end{aligned}$$

By definition of $\tilde{\xi}$

$$\tilde{\xi}_{\tilde{S}} = \ell_j M_{\tilde{S}j}, \quad \tilde{\xi}_j = M_{j\tilde{S}} \ell_{\tilde{S}} + \ell_j M_{jj} - \ell_j = -\xi_j - \ell_j, \quad \tilde{\xi}_{S^c} = \ell_j M_{S^c j},$$

and thus,

$$\tilde{\xi} = \ell_j \beta - \xi_j e_j.$$

This implies that

$$\xi^- = \xi - \xi_j e_j + \frac{\beta_{\tilde{S}}^T \ell_{\tilde{S}} + \ell_j M_{jj}}{M_{jj}} \beta = \mathcal{R}_j(\xi) - \frac{\xi_j}{M_{jj}} \beta.$$

For D_ℓ^- , we have

$$D_\ell^- = \mathbf{1}_{\tilde{S}}^T \xi_{\tilde{S}}^- = D_\ell - \xi_j - \frac{\xi_j}{M_{jj}} \cdot \mathbf{1}_{\tilde{S}}^T \beta_{\tilde{S}} = D_\ell - \frac{\xi_j (\mathbf{1}_{\tilde{S}}^T \beta_{\tilde{S}} + M_{jj})}{M_{jj}} = D_\ell - \frac{\xi_j \tilde{\eta}_j}{M_{jj}}.$$

■

The implementation of UPDATE_TILDE_SHRINK_SUPPORT is summarized in Algorithm 14.

Algorithm 14 UPDATE_UTILDE_SHRINK_SUPPORT

Inputs: Original support S , new index j , vector ℓ , intermediate variables $\text{Par}_1, \text{Par}_3$.

Procedure:

- 1: $\beta_{\tilde{S}} \leftarrow M_{j\tilde{S}}^T, \beta_{\tilde{S}^c} \leftarrow \begin{pmatrix} -1 \\ M_{\tilde{S}^c j} \end{pmatrix}, \tilde{\beta}_{\tilde{S}} \leftarrow M_{j\tilde{S}}^T, \tilde{\beta}_{\tilde{S}^c} \leftarrow 0$;
- 2: $D \leftarrow D - \frac{\tilde{\eta}_j^2}{M_{jj}}$;
- 3: $D_\ell \leftarrow D_\ell - \frac{\xi_j \tilde{\eta}_j}{M_{jj}}$;
- 4: $\xi \leftarrow \mathcal{R}_j(\xi) - \frac{\xi_j}{M_{jj}} \beta$;
- 5: $\tilde{\eta} \leftarrow \mathcal{R}_j(\tilde{\eta}) - \frac{\tilde{\eta}_j}{M_{jj}} \beta$;
- 6: $M_{\cdot, \tilde{S}} \leftarrow \mathcal{R}_j(M_{\cdot, \tilde{S}}) - \frac{1}{M_{jj}} \beta \tilde{\beta}_{\tilde{S}}^T, M_{\cdot, j} \leftarrow 0$.

Output: $\text{Par}_1, \text{Par}_3$.

Algorithm 15 DIRECT_UTILDE_UPDATE

Inputs: Support S , vector-update-vector ℓ , intermediate variables M .

Procedure:

- 1: $\xi_S \leftarrow -M_{SS}\ell_S$, $\xi_{S^c} \leftarrow -\ell_{S^c} - M_{S^cS}\ell_S$;
- 2: $D_\ell \leftarrow \mathbf{1}_S^T \xi_S$.

Output: ξ, D_ℓ .

C.4.4 DIRECT_UTILDE_UPDATE

At the beginning of each time t , we need to recompute ξ and D_ℓ . The implementation is summarized in Algorithm 15.

C.5 Complexity Analysis

Similar to Appendix B, we can analyze the computation complexity. The analysis here is much simpler than the last case since the implementation is quite straightforward. Table 6 summarizes the results.

Table 6: Computation complexity of each sub-routine in Algorithm 10.

	(Wz) -type	$(z\tilde{z}^T)$ -type	$(z^T z)$ -type	(az) -type
FIND_UTILDE_LAMBDA	0	0	$2n$	0
UPDATE_BY_UTILDE_LAMBDA	0	0	0	n
UPDATE_UTILDE_EXPAND_SUPPORT	$(n-s-1)s$	$n(s+1)$	s	$2n$
UPDATE_UTILDE_SHRINK_SUPPORT	0	$n(s-1)$	0	$2n$
DIRECT_UTILDE_UPDATE	ns	0	s	0

Let k_r^+ and k_r^- be the number of turning points that S is expanded and shrunked respectively and $k_r = k_r^+ + k_r^-$ be the total number of tuning points, then the complexity is

$$C_{2t} \leq ns + n + 3nk_r + n(2s+3)k_r^+ + n(s+1)k_r^- + O(k_r) = ns(2k_r+1) + n(6k_r+1) + O(k_r).$$

D Implementation of HOP Algorithm With Time-Varying A, r

D.1 Intermediate Variables

Based on the results in Appendix B and Appendix C, we can concatenate Algorithm 3 and Algorithm 10. Thus we define $\text{Par}_1, \text{Par}_2, \text{Par}_3$ as $\text{Par}_1 = \{M, \tilde{\eta}, D\}$, $\text{Par}_2 = \{\eta, D_g, D_{gg}, D_{gr}\}$, $\text{Par}_3 = \{\xi, D_\ell\}$ where all parameters are defined in previous appendices.

D.2 Implementation

Note that only two sub-routines involves the matrix A , namely UPDATE_EXPAND_SUPPORT and UPDATE_UTILDE_EXPAND_SUPPORT, and moreover they only involve the j -th column of A . Thus, we can use the late update of A as in Algorithm 9 for acceleration. Algorithm 16 below describes the implementation.

D.3 Complexity Analysis

The complexity of Algorithm 16 is just the sum of that of Algorithm 3 and Algorithm 10, i.e.

$$C_t = C_{1t} + C_{2t} = ns_* + ns(3k_A + 2k_r + 2) + n(12k_A + 6k_r + 3) + O(k_A + k_r).$$

Algorithm 16 HOP Algorithm for time-varying A, r with late update of A

Inputs: Initial parameters $A^{(0)}$, vectors $\{r^{(t)} : t = 1, 2, \dots\}$,
matrix-update-vectors $\{g^{(t)}, t = 1, 2, \dots\}$.

Initialization:

$x \leftarrow$ as the optimum corresponding to $A^{(0)}, r^{(0)}$.
 $S \leftarrow \text{supp}(x), S_* \leftarrow S$;
Calculate (x, μ, μ_0) via (10)-(12)
 $v \leftarrow (x_S, -\mu_{S^c})$;
Calculate intermediate variables $(\text{Par}_1, \text{Par}_2)$ via (18)-(20) based on $r^{(0)}, g^{(1)}$;

Procedure:

```
1: for  $t = 1, 2, \dots$  do
2:    $\lambda \leftarrow 0$ ;
3:   while  $\lambda < 1$  do
4:      $(\lambda^{\text{inc}}, j, S^{\text{new}}) \leftarrow \text{FIND\_LAMBDA}(S, v; \text{Par}_1, \text{Par}_2)$ ;
5:      $\lambda^{\text{inc}} \leftarrow \min\{\lambda^{\text{inc}}, 1 - \lambda\}$ ;
6:      $\lambda \leftarrow \lambda + \lambda^{\text{inc}}$ ;
7:      $(v, \mu_0; \text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_BY\_LAMBDA}(\lambda^{\text{inc}}; v, \mu_0, \text{Par}_1, \text{Par}_2)$ ;
8:     if  $S^{\text{new}} = S \cup \{j\}$  then
9:        $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_EXPAND\_SUPPORT}(\lambda, S, j; A, r^{(t-1)}, g^{(t)}, \text{Par}_1, \text{Par}_2)$ ;
10:      if  $j \notin S_*$  then
11:         $G \leftarrow (g^{(1)}, \dots, g^{(t-1)})$ ;
12:         $A_{\cdot, j} \leftarrow A_{\cdot, j} + GG_{j, \cdot}^T$ ;
13:         $S_* = S_* \cup \{j\}$ ;
14:      end if
15:    else if  $S^{\text{new}} = S \setminus \{j\}$  then
16:       $(\text{Par}_1, \text{Par}_2) \leftarrow \text{UPDATE\_SHRINK\_SUPPORT}(S, j; r^{(t-1)}, g^{(t)}, \text{Par}_1, \text{Par}_2)$ ;
17:    end if
18:     $S \leftarrow S^{\text{new}}$ ;
19:  end while
20:   $A_{\cdot, S_*} \leftarrow A_{\cdot, S_*} + g^{(t)}(g_{S_*}^{(t)})^T$ ;
21:   $\ell^{(t)} \leftarrow r^{(t)} - r^{(t-1)}$ ;
22:   $\text{Par}_3 \leftarrow \text{DIRECT\_UTILDE\_UPDATE}(S, \text{Par}_1, \ell^{(t)})$ ;
23:   $\underline{\lambda} \leftarrow 0$ ;
24:  while  $\underline{\lambda} < 1$  do
25:     $(\underline{\lambda}^{\text{inc}}, j, S^{\text{new}}) \leftarrow \text{FIND\_UTILDE\_LAMBDA}(v; \text{Par}_1, \text{Par}_3)$ ;
26:     $\underline{\lambda}^{\text{inc}} \leftarrow \min\{\underline{\lambda}^{\text{inc}}, 1 - \underline{\lambda}\}$ ;
27:     $(v, \mu_0) \leftarrow \text{UPDATE\_BY\_UTILDE\_LAMBDA}(\underline{\lambda}^{\text{inc}}; v, \mu_0, \text{Par}_1, \text{Par}_3)$ ;
28:    if  $S^{\text{new}} = S \cup \{j\}$  then
29:       $(\text{Par}_1, \text{Par}_3) \leftarrow \text{UPDATE\_UTILDE\_EXPAND\_SUPPORT}(S, j, A, \ell^{(t)}; \text{Par}_1, \text{Par}_3)$ ;
30:      if  $j \notin S_*$  then
31:         $G \leftarrow (g^{(1)}, \dots, g^{(t-1)})$ ;
32:         $A_{\cdot, j} \leftarrow A_{\cdot, j} + GG_{j, \cdot}^T$ ;
33:         $S_* = S_* \cup \{j\}$ ;
34:      end if
35:    else if  $S^{\text{new}} = S \setminus \{j\}$  then
36:       $(\text{Par}_1, \text{Par}_3) \leftarrow \text{UPDATE\_UTILDE\_SHRINK\_SUPPORT}(S, j, \ell^{(t)}; \text{Par}_1, \text{Par}_3)$ ;
37:    end if
38:     $S \leftarrow S^{\text{new}}$ ;
39:     $\underline{\lambda} \leftarrow \underline{\lambda} + \underline{\lambda}^{\text{inc}}$ .
40:  end while
41:   $\text{Par}_2 \leftarrow \text{DIRECT\_UPDATE}(S, r^{(t)}, g^{(t+1)}; \text{Par}_1, \text{Par}_2)$ ;
42:   $x_S^{(t)} \leftarrow x_S, \quad x_{S^c}^{(t)} \leftarrow 0$ .
43: end for
```

Output: $x^{(1)}, x^{(2)}, \dots$
